

Kernels & Probabilistic Modeling

Mengye Ren

(Slides credit to David Rosenberg, He He, et al.)

NYU

October 1, 2024



- Today (Oct 1): Kernels and Probabilistic Modeling
- Oct 8: Guest Lecture
- Oct 15: Homework 2 Due
- Oct 15: Legislative Day No Class
- Oct 22: Midterm, in class, closed-book, covers everything including Oct 8

HW1 + HW2

Expressivity of Hypothesis Space

$$\hat{y} = \underbrace{w^T}_{\text{weights}} \underbrace{x}_{\text{features}} = \phi(x)$$

- For linear models, to grow the hypothesis spaces, we must add features.
- Sometimes we say a larger hypothesis is **more expressive**.
 - (can fit more relationships between input and action)
- Many ways to create new features.

Handling Nonlinearity with Linear Methods

Example Task: Predicting Health

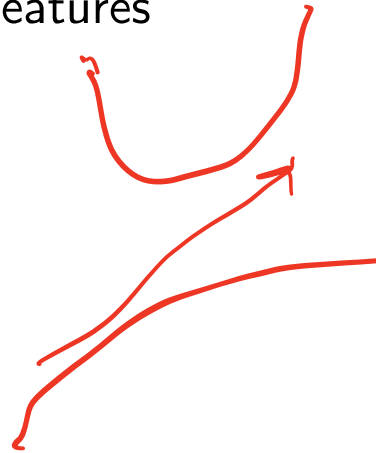
- General Philosophy: Extract every feature that might be relevant
- Features for medical diagnosis
 - height ✓
 - weight ✓
 - body temperature ✓
 - blood pressure
 - etc...

Feature Issues for Linear Predictors

- For linear predictors, it's important **how** features are added
 - The relation between a feature and the label may not be linear
 - There may be complex dependence among features

Feature Issues for Linear Predictors

- For linear predictors, it's important **how** features are added
 - The relation between a feature and the label may not be linear
 - There may be complex dependence among features
- Three types of nonlinearities can cause problems:
 - Non-monotonicity
 - Saturation
 - Interactions between features



Non-monotonicity: The Issue

- Feature Map: $\phi(x) = [1, \text{temperature}(x)]$
- Action: Predict health score $y \in \mathbb{R}$ (positive is good)
- Hypothesis Space $\mathcal{F} = \{\text{affine functions of temperature}\}$

$$\begin{aligned} w^T x + b \\ = \\ w^T [x, 1] \end{aligned}$$

Non-monotonicity: The Issue

- Feature Map: $\phi(x) = [1, \text{temperature}(x)]$
- Action: Predict health score $y \in \mathbb{R}$ (positive is good)
- Hypothesis Space $\mathcal{F} = \{\text{affine functions of temperature}\}$
- Issue:
 - Health is not an affine function of temperature.

Non-monotonicity: The Issue

- Feature Map: $\phi(x) = [1, \text{temperature}(x)]$
- Action: Predict health score $y \in \mathbb{R}$ (positive is good)
- Hypothesis Space $\mathcal{F} = \{\text{affine functions of temperature}\}$
- Issue:
 - Health is not an affine function of temperature.
 - Affine function can either say
 - Very high is bad and very low is good, or
 - Very low is bad and very high is good,
 - But here, both extremes are bad.

Non-monotonicity: Solution 1

- Transform the input:

$$\phi(x) = \left[1, \underbrace{\{\text{temperature}(x) - 37\}^2} \right],$$

where 37 is “normal” temperature in Celsius.

Non-monotonicity: Solution 1

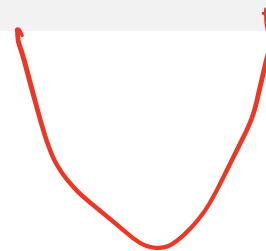
- Transform the input:

$$\phi(x) = \left[1, \{\text{temperature}(x) - 37\}^2 \right],$$

where 37 is “normal” temperature in Celsius.

- Ok, but requires manually-specified domain knowledge
 - Do we really need that?
 - What does $w^T \phi(x)$ look like?

Non-monotonicity: Solution 2



- Think less, put in more:

$$\phi(x) = \left[\underline{1, \text{temperature}(x), \{\text{temperature}(x)\}^2} \right].$$

- More expressive than Solution 1.

General Rule

Features should be simple building blocks that can be pieced together.

Saturation: The Issue

- Setting: Find products relevant to user's query

Saturation: The Issue

- Setting: Find products relevant to user's query
- Input: Product x
- Output: Score the relevance of x to user's query

Saturation: The Issue

- Setting: Find products relevant to user's query
- Input: Product x
- Output: Score the relevance of x to user's query
- Feature Map:

$$\phi(x) = [1, N(x)],$$

where $N(x)$ = number of people who bought x .

Saturation: The Issue

- Setting: Find products relevant to user's query
- Input: Product x
- Output: Score the relevance of x to user's query
- Feature Map:

$$\phi(x) = [1, \underline{N(x)}], \quad \text{---} \log N(x)$$

where $N(x)$ = number of people who bought x .

- We expect a monotonic relationship between $N(x)$ and relevance, but also expect diminishing return.

Saturation: Solve with nonlinear transform

- Smooth nonlinear transformation:

$$\phi(x) = [1, \log\{1 + N(x)\}]$$

- $\log(\cdot)$ good for values with large dynamic ranges

Saturation: Solve with nonlinear transform

- Smooth nonlinear transformation:

$$\phi(x) = [1, \log\{1 + N(x)\}]$$

- $\log(\cdot)$ good for values with large dynamic ranges
- Discretization (a discontinuous transformation):

$$\phi(x) = (\mathbb{1}[0 \leq N(x) < 10], \mathbb{1}[10 \leq N(x) < 100], \dots)$$

- Small buckets allow quite flexible relationship

Interactions: The Issue

- Input: Patient information x
- Action: Health score $y \in \mathbb{R}$ (higher is better)
- Feature Map

$$\phi(x) = [\text{height}(x), \text{weight}(x)]$$

The diagram illustrates the feature map $\phi(x) = [\text{height}(x), \text{weight}(x)]$. Red arrows point from $\text{height}(x)$ to w_1 and from $\text{weight}(x)$ to w_2 . A red box encloses w_1 and w_2 .

Interactions: The Issue

- Input: Patient information x
- Action: Health score $y \in \mathbb{R}$ (higher is better)
- Feature Map

$$\phi(x) = [\text{height}(x), \text{weight}(x)]$$

- Issue: It's the weight *relative* to the height that's important.

Interactions: The Issue

- Input: Patient information x
- Action: Health score $y \in \mathbb{R}$ (higher is better)
- Feature Map

$$\phi(x) = [\text{height}(x), \text{weight}(x)]$$

- Issue: It's the weight *relative* to the height that's important.
- Impossible to get with these features and a linear classifier.
- Need some **interaction** between height and weight.

Interactions: Approach 1

- Google “ideal weight from height”
- J. D. Robinson’s “ideal weight” formula:

$$\text{weight}(\text{kg}) = 52 + 1.9 [\text{height}(\text{in}) - 60]$$

Interactions: Approach 1

- Google “ideal weight from height”
- J. D. Robinson’s “ideal weight” formula:

$$\text{weight}(\text{kg}) = 52 + 1.9 [\text{height}(\text{in}) - 60]$$

- Make score square deviation between height(h) and ideal weight(w)

$$f(x) = (52 + 1.9 [h(x) - 60] - w(x))^2$$

Interactions: Approach 1

- Google “ideal weight from height”
- J. D. Robinson’s “ideal weight” formula:

$$\text{weight}(\text{kg}) = 52 + 1.9 [\text{height}(\text{in}) - 60]$$

- Make score square deviation between height(h) and ideal weight(w)

$$f(x) = (52 + 1.9 [h(x) - 60] - w(x))^2$$

- WolframAlpha for complicated Mathematics:

$$f(x) = 3.61h(x)^2 - 3.8h(x)w(x) - 235.6h(x) + w(x)^2 + 124w(x) + 3844$$

Interactions: Approach 2

- Just include all second order features:

$$\phi(x) = \left[1, h(x), w(x), h(x)^2, w(x)^2, \underbrace{h(x)w(x)}_{\text{cross term}} \right]$$

- More flexible, no Google, no WolframAlpha.

General Principle

Simpler building blocks replace a single “smart” feature.

Monomial Interaction Terms

Interaction terms are useful building blocks to model non-linearities in features.

- Suppose we start with $x = (1, x_1, \dots, x_d) \in \mathbb{R}^{d+1} = \mathcal{X}$.

Monomial Interaction Terms

Interaction terms are useful building blocks to model non-linearities in features.

- Suppose we start with $x = (1, x_1, \dots, x_d) \in \mathbb{R}^{d+1} = \mathcal{X}$.
- Consider adding all **monomials** of degree M : $x_1^{p_1} \dots x_d^{p_d}$, with $p_1 + \dots + p_d = M$.
 - Monomials with degree 2 in 2D space: x_1^2, x_2^2, x_1x_2

Monomial Interaction Terms

Interaction terms are useful building blocks to model non-linearities in features.

- Suppose we start with $x = (1, x_1, \dots, x_d) \in \mathbb{R}^{d+1} = \mathcal{X}$.
- Consider adding all **monomials** of degree M : $x_1^{p_1} \dots x_d^{p_d}$, with $p_1 + \dots + p_d = M$.
 - Monomials with degree 2 in 2D space: x_1^2, x_2^2, x_1x_2
- How many features will we end up with?

d features to degree M .
3 feature 2. $\binom{n}{2}$

x_1x_2 x_2x_3 x_1^2 x_2^2 x_3^2 x_1x_3 $d!$

Monomial Interaction Terms

Interaction terms are useful building blocks to model non-linearities in features.

- Suppose we start with $x = (1, x_1, \dots, x_d) \in \mathbb{R}^{d+1} = \mathcal{X}$.
- Consider adding all **monomials** of degree M : $x_1^{p_1} \dots x_d^{p_d}$, with $p_1 + \dots + p_d = M$.
 - Monomials with degree 2 in 2D space: x_1^2, x_2^2, x_1x_2
- How many features will we end up with?

$$\begin{array}{l} x_1^3 \leftrightarrow * \cdot * \cdot * \mid \mid \\ x_1^2 x_2 \leftrightarrow * \cdot * \mid * \mid \\ x_1^2 x_3 \leftrightarrow * \cdot * \mid \mid * \\ x_1 x_2^2 \leftrightarrow * \mid * \cdot * \mid \\ \boxed{x_1 x_2 x_3} \leftrightarrow * \mid * \mid * \\ x_1 x_3^2 \leftrightarrow * \mid \mid * \cdot * \\ x_2^3 \leftrightarrow \mid * \cdot * \cdot * \mid \\ x_2^2 x_3 \leftrightarrow \mid * \cdot * \mid * \\ x_2 x_3^2 \leftrightarrow \mid * \mid * \cdot * \\ x_3^3 \leftrightarrow \mid \mid * \cdot * \cdot * \end{array}$$

Monomial Interaction Terms

Interaction terms are useful building blocks to model non-linearities in features.

- Suppose we start with $x = (1, x_1, \dots, x_d) \in \mathbb{R}^{d+1} = \mathcal{X}$.
- Consider adding all **monomials** of degree M : $x_1^{p_1} \dots x_d^{p_d}$, with $p_1 + \dots + p_d = M$.
 - Monomials with degree 2 in 2D space: x_1^2, x_2^2, x_1x_2
- How many features will we end up with?

$$\begin{aligned}x_1^3 &\leftrightarrow * * * | | \\x_1^2x_2 &\leftrightarrow * * | * | \\x_1^2x_3 &\leftrightarrow * * | | * \\x_1x_2^2 &\leftrightarrow * | * * | \\x_1x_2x_3 &\leftrightarrow * | * | * \\x_1x_3^2 &\leftrightarrow * | | * * \\x_2^3 &\leftrightarrow | * * * | \\x_2^2x_3 &\leftrightarrow | * * | * \\x_2x_3^2 &\leftrightarrow | * | * * \\x_3^3 &\leftrightarrow | | * * *\end{aligned}$$

Big Feature Spaces

This leads to extremely **large data matrices**

- For $d = 40$ and $M = 8$, we get 314457495 features.

Big Feature Spaces

This leads to extremely **large data matrices**

- For $d = 40$ and $M = 8$, we get 314457495 features.

Very large feature spaces have two potential issues:

- Overfitting
- Memory and computational costs

Big Feature Spaces

This leads to extremely **large data matrices**

- For $d = 40$ and $M = 8$, we get 314457495 features.

Very large feature spaces have two potential issues:

- Overfitting \longrightarrow *# parameter = # features.*
- Memory and computational costs

Solutions:

- Overfitting we handle with regularization.
- Kernel methods can help with memory and computational costs when we go to high (or infinite) dimensional spaces. $\gg \gg \gg >$ 2fb8a69 (05)

The Kernel Trick

SVM with Explicit Feature Map

- Let $\psi : \mathcal{X} \rightarrow \mathbb{R}^d$ be a feature map.
- The SVM objective (with explicit feature map):

$$\min_{w \in \mathbb{R}^d} \underbrace{\frac{1}{2} \|w\|^2}_{L_2 \text{ reg}} + \frac{c}{n} \sum_{i=1}^n \underbrace{\max(0, 1 - y_i w^T \psi(x_i))}_{\text{hinge loss}}.$$

SVM with Explicit Feature Map

- Let $\psi : \mathcal{X} \rightarrow \mathbb{R}^d$ be a feature map.
- The SVM objective (with explicit feature map):

$$\min_{w \in \mathbb{R}^d} \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \max(0, 1 - y_i w^T \psi(x_i)).$$

Handwritten notes: $\psi(x)$ with an arrow pointing to the $\psi(x_i)$ term, a downward arrow, and a bracket containing $[1, x_1 x_2 \dots]$.

- Computation is costly if d is large (e.g. with high-degree monomials)

SVM with Explicit Feature Map

- Let $\psi : \mathcal{X} \rightarrow \mathbb{R}^d$ be a feature map.
- The SVM objective (with explicit feature map):

$$\min_{w \in \mathbb{R}^d} \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \max(0, 1 - y_i w^T \psi(x_i)).$$

- Computation is costly if d is large (e.g. with high-degree monomials)
- Last time we mentioned an equivalent optimization problem from Lagrangian duality.

SVM Dual Problem

$$\max \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \phi(x_j)^T \phi(x_i)$$

• By Lagrangian duality, it is equivalent to solve the following dual problem:

• If α^* is an optimal value, then

$$w^* = \sum_i \alpha_i^* y_i \phi(x_i)$$

$$\hat{y} = \sum_i \alpha_i^* y_i \phi(x_i) \phi(x)$$

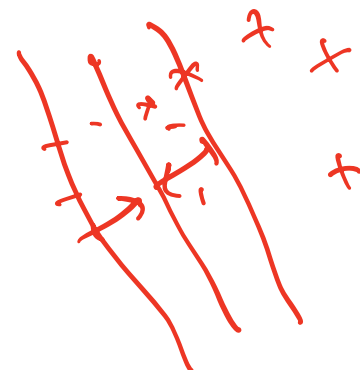
w^* test input

s.t.

$$\sum_i \alpha_i y_i = 0$$

$$\alpha_i \in [0, \frac{1}{n}]$$

class balancing.



SVM Dual Problem

Learning: $\max \sum \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \psi(x_i)^T \psi(x_j)$

• By Lagrangian duality, it is equivalent to solve the following dual problem:

• If α^* is an optimal value, then

$d \gg \eta$ $(d \times n)$

$(d \times n)^T (d \times n)$
 $\begin{matrix} d \\ \hline n \times n \end{matrix}$

• **Key observation:** $\psi(x)$ only shows up in **inner products** with another $\psi(x')$ for both training and inference.

pred: $\hat{y}(x) = \sum_i \alpha_i^* y_i \psi(x_i)^T \psi(x)$

inference.

train ex test ex

Compute the Inner Products

$$\psi(x) \mapsto \mathbb{R}^d \times d.$$

Consider 2D data. Let's introduce **degree-2 monomials** using $\psi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$.

$$(x_1, x_2) \mapsto (x_1^2, \sqrt{2}x_1x_2, x_2^2).$$

The inner product is

We can calculate the inner product $\psi(x)^T \psi(x')$ in the original input space without accessing the features $\psi(x)$!

$$\begin{aligned} \psi(x)^T \psi(x') &= x_1^2 x_1'^2 + (\sqrt{2} x_1 x_2) (\sqrt{2} x_1' x_2') + x_2^2 x_2'^2 \\ &= x_1^2 x_1'^2 + 2 x_1 x_1' x_2 x_2' + x_2^2 x_2'^2 \\ &= (x_1 x_1' + x_2 x_2')^2 \\ &= (x^T x')^2 \rightarrow n \times n \end{aligned}$$

Compute the Inner Products

Now, consider **monomials up to degree-2**:

$$(x_1, x_2) \mapsto (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2).$$

The inner product can be computed by

$$\psi(x)^T \psi(x') = (1 + x^T x')^2 \text{ (check).}$$

$$k(x, x') = \textcircled{S}$$

Compute the Inner Products

Now, consider **monomials up to degree-2**:

$$(x_1, x_2) \mapsto (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2).$$

The inner product can be computed by

$$\psi(x)^T \psi(x') = (1 + x^T x')^2 \quad (\text{check}).$$

More generally, for features maps producing monomials up to degree- p , we have

$$\psi(x)^T \psi(x') = (1 + x^T x')^p.$$

(Note that the coefficients of each monomial in ψ may not be 1)

Kernel Trick

$$\odot \psi(x)$$

$$\psi(x)^T \psi(x')$$

Kernel trick: we do not need explicit features to calculate inner products.

- Using explicit features: $O(d^p)$
- Using implicit computation: $O(d) \left(\frac{K(x, x')^p}{d} \right)$

Kernel Function

The Kernel Function

- Input space: \mathcal{X}
- Feature space: \mathcal{H}
- Feature map: $\psi : \mathcal{X} \rightarrow \mathcal{H}$
- The kernel function corresponding to ψ is

$$k(x, x') = \langle \psi(x), \psi(x') \rangle,$$

where $\langle \cdot, \cdot \rangle$ is the inner product associated with \mathcal{H} .

The Kernel Function

- **Input space:** \mathcal{X}
- **Feature space:** \mathcal{H}
- **Feature map:** $\psi : \mathcal{X} \rightarrow \mathcal{H}$
- The **kernel function** corresponding to ψ is

$$\underline{k(x, x')} = \langle \psi(x), \psi(x') \rangle,$$

where $\langle \cdot, \cdot \rangle$ is the inner product associated with \mathcal{H} .

Why introduce this new notation $k(x, x')$?

The Kernel Function

- **Input space:** \mathcal{X}
- **Feature space:** \mathcal{H}
- **Feature map:** $\psi : \mathcal{X} \rightarrow \mathcal{H}$
- The **kernel function** corresponding to ψ is

$$k(x, x') = \langle \psi(x), \psi(x') \rangle,$$

where $\langle \cdot, \cdot \rangle$ is the inner product associated with \mathcal{H} .

Why introduce this new notation $k(x, x')$?

- We can often evaluate $k(x, x')$ without explicitly computing $\psi(x)$ and $\psi(x')$.

The Kernel Function

- **Input space:** \mathcal{X}
- **Feature space:** \mathcal{H}
- **Feature map:** $\psi : \mathcal{X} \rightarrow \mathcal{H}$
- The **kernel function** corresponding to ψ is

$$k(x, x') = \langle \psi(x), \psi(x') \rangle,$$

where $\langle \cdot, \cdot \rangle$ is the inner product associated with \mathcal{H} .

Why introduce this new notation $k(x, x')$?

- We can often evaluate $k(x, x')$ without explicitly computing $\psi(x)$ and $\psi(x')$.

When can we use the kernel trick?

Some Methods Can Be “Kernelized”

Definition

A method is kernelized if every feature vector $\psi(x)$ only appears inside an inner product with another feature vector $\psi(x')$. This applies to both the optimization problem and the prediction function.

learning

inference

Some Methods Can Be “Kernelized”

Definition

A method is **kernelized** if every feature vector $\psi(x)$ only appears inside an inner product with another feature vector $\psi(x')$. This applies to both the optimization problem and the prediction function.

The SVM Dual is a kernelization of the original SVM formulation.

Optimization:

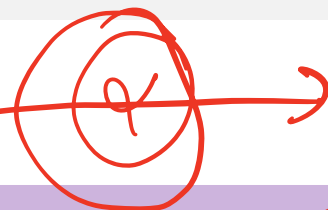
$$\begin{aligned} \text{maximize} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \psi(x_j)^T \psi(x_i) \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \quad \text{and} \quad \alpha_i \in \left[0, \frac{c}{n}\right] \quad \forall i. \end{aligned}$$

Prediction:

$$\hat{f}(x) = \sum_{i=1}^n \alpha_i^* y_i \psi(x_i)^T \psi(x).$$

The Kernel Matrix

$$[1 \oplus h^2 \oplus w^2 \oplus hw]$$



$$(1 \oplus x^T x')^2$$

(h, w)

Definition

The **kernel matrix** for a kernel k on $x_1, \dots, x_n \in \mathcal{X}$ is

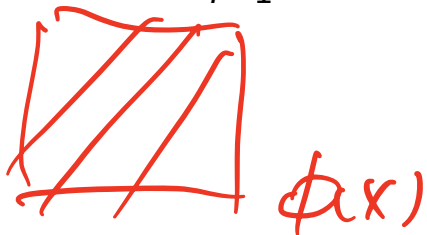
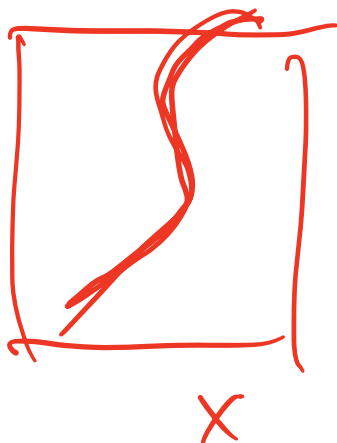
$$K = (k(x_i, x_j))_{i,j} = \begin{pmatrix} k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & \ddots & \cdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) \end{pmatrix} \in \mathbb{R}^{n \times n}.$$

- In ML this is also called a Gram matrix, but traditionally (in linear algebra), Gram matrices are defined without reference to a kernel or feature map.

The Kernel Matrix

- The kernel matrix **summarizes all the information** we need about the training inputs x_1, \dots, x_n to solve a kernelized optimization problem.
- In the kernelized SVM, we can replace $\psi(x_i)^T \psi(x_j)$ with K_{ij} :

$$\begin{aligned} \text{maximize}_{\alpha} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j K_{ij} \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \quad \text{and} \quad \alpha_i \in \left[0, \frac{c}{n}\right] \quad i = 1, \dots, n. \end{aligned}$$



$$x_i^T x_j$$

$$\phi(x_i)^T \phi(x_j)$$

Gram Matrix

$$\phi(x)^T \phi(x)$$

Kernel Methods

Given a kernelized ML algorithm (i.e. all $\psi(x)$'s show up as $\langle \psi(x), \psi(x') \rangle$),

- Can swap out the inner product for a new kernel function.
- New kernel may correspond to a **very high-dimensional** feature space.

Kernel Methods

Given a kernelized ML algorithm (i.e. all $\psi(x)$'s show up as $\langle \psi(x), \psi(x') \rangle$),

- Can swap out the inner product for a new kernel function.
- New kernel may correspond to a **very high-dimensional** feature space.
- Once the kernel matrix is computed, the computational cost **depends on number of data points** n , rather than the dimension of feature space d .
- Useful when $d \gg n$.

Kernel Methods

Given a kernelized ML algorithm (i.e. all $\psi(x)$'s show up as $\langle \psi(x), \psi(x') \rangle$),

- Can swap out the inner product for a new kernel function.
- New kernel may correspond to a **very high-dimensional** feature space.
- Once the kernel matrix is computed, the computational cost **depends on number of data points** n , rather than the dimension of feature space d .
- Useful when $d \gg n$.
- Computing the kernel matrix may still depend on d and the essence of the trick is getting around this $O(d)$ dependence.

Example Kernels

Kernels as Similarity Scores

- Often useful to think of the $k(x, x')$ as a **similarity score** for x and x' .
- We can design similarity functions without thinking about the explicit feature map, e.g. “string kernels”, “graph kernels”.
- How do we know that our kernel functions actually correspond to inner products in some feature space?

How to Get Kernels?

- Explicitly construct $\psi(x) : \mathcal{X} \rightarrow \mathbb{R}^d$ (e.g. monomials) and define $k(x, x') = \psi(x)^T \psi(x')$.
- Directly define the kernel function $k(x, x')$ (“similarity score”), and **verify it corresponds to $\langle \psi(x), \psi(x') \rangle$ for some ψ .**

There are many theorems to help us with the second approach.

Linear Algebra Review: Positive Semidefinite Matrices

Definition

A real, symmetric matrix $M \in \mathbb{R}^{n \times n}$ is **positive semidefinite (psd)** if for any $x \in \mathbb{R}^n$,

$$x^T M x \geq 0.$$

Theorem

The following conditions are each necessary and sufficient for a symmetric matrix M to be positive semidefinite:

- M can be factorized as $M = R^T R$ for some matrix R .
- All eigenvalues of M are greater than or equal to 0.

Positive Definite Kernel

Definition

A symmetric function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a **positive definite (pd)** kernel on \mathcal{X} if for any finite set $\{x_1, \dots, x_n\} \in \mathcal{X}$ ($n \in \mathbb{N}$), the kernel matrix on this set

$$K = (k(x_i, x_j))_{i,j} = \begin{pmatrix} k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & \ddots & \cdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) \end{pmatrix}$$

is a positive semidefinite matrix.

- Symmetric: $k(x, x') = k(x', x)$
- The kernel matrix needs to be positive semidefinite for any finite set of points.
- Equivalent definition: $\sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(x_i, x_j) \geq 0$ given $\alpha_i \in \mathbb{R} \forall i$.

Mercer's Theorem

Theorem

A symmetric function $k(x, x')$ can be expressed as an inner product

$$k(x, x') = \langle \psi(x), \psi(x') \rangle$$

for some ψ if and only if $k(x, x')$ is **positive definite**.

- Proving a kernel function is positive definite is typically not easy.
- But we can construct new kernels from valid kernels.

Generating New Kernels from Old

- Suppose $k, k_1, k_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ are pd kernels. Then so are the following:

$$k_{\text{new}}(x, x') = \alpha k(x, x') \quad \text{for } \underline{\alpha \geq 0} \quad (\text{non-negative scaling})$$

Generating New Kernels from Old

- Suppose $k, k_1, k_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ are pd kernels. Then so are the following:

$$k_{\text{new}}(x, x') = \alpha k(x, x') \quad \text{for } \alpha \geq 0 \quad (\text{non-negative scaling})$$

$$k_{\text{new}}(x, x') = \underline{k_1(x, x')} + \underline{k_2(x, x')} \quad (\text{sum})$$

Generating New Kernels from Old

- Suppose $k, k_1, k_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ are pd kernels. Then so are the following:

$$k_{\text{new}}(x, x') = \alpha k(x, x') \quad \text{for } \alpha \geq 0 \quad (\text{non-negative scaling})$$

$$k_{\text{new}}(x, x') = k_1(x, x') + k_2(x, x') \quad (\text{sum})$$

$$k_{\text{new}}(x, x') = k_1(x, x') k_2(x, x') \quad (\text{product})$$

Generating New Kernels from Old

- Suppose $k, k_1, k_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ are pd kernels. Then so are the following:

$$k_{\text{new}}(x, x') = \alpha k(x, x') \quad \text{for } \alpha \geq 0 \quad (\text{non-negative scaling})$$

$$k_{\text{new}}(x, x') = k_1(x, x') + k_2(x, x') \quad (\text{sum})$$

$$k_{\text{new}}(x, x') = k_1(x, x') k_2(x, x') \quad (\text{product})$$

$$k_{\text{new}}(x, x') = k(\underbrace{\psi(x)}, \underbrace{\psi(x')}) \quad \text{for any function } \psi(\cdot) \quad (\text{recursion})$$

Generating New Kernels from Old

- Suppose $k, k_1, k_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ are pd kernels. Then so are the following:

$$k_{\text{new}}(x, x') = \alpha k(x, x') \quad \text{for } \alpha \geq 0 \quad (\text{non-negative scaling})$$

$$k_{\text{new}}(x, x') = k_1(x, x') + k_2(x, x') \quad (\text{sum})$$

$$k_{\text{new}}(x, x') = k_1(x, x') k_2(x, x') \quad (\text{product})$$

$$k_{\text{new}}(x, x') = k(\psi(x), \psi(x')) \quad \text{for any function } \psi(\cdot) \quad (\text{recursion})$$

$$k_{\text{new}}(x, x') = \underbrace{f(x)} \underbrace{f(x')} \quad \text{for any function } f(\cdot) \quad (f \text{ as 1D feature map})$$

- Lots more theorems to help you construct new kernels from old.

Linear Kernel

- Input space: $\mathcal{X} = \mathbb{R}^d$
- Feature space: $\mathcal{H} = \mathbb{R}^d$, with standard inner product
- Feature map

$$\psi(x) = x$$

- Kernel:

$$k(x, x') = \underline{x^T x'}$$

Quadratic Kernel in \mathbb{R}^d

- Input space $\mathcal{X} = \mathbb{R}^d$
- Feature space: $\mathcal{H} = \mathbb{R}^D$, where $D = d + \binom{d}{2} \approx d^2/2$.
- Feature map:

$$\psi(x) = (x_1, \dots, x_d, x_1^2, \dots, x_d^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_ix_j, \dots, \sqrt{2}x_{d-1}x_d)^T$$

- Then for $\forall x, x' \in \mathbb{R}^d$

$$\begin{aligned} k(x, x') &= \langle \psi(x), \psi(x') \rangle \\ &= \langle x, x' \rangle + \langle x, x' \rangle^2 \end{aligned}$$

- Computation for inner product with explicit mapping: $O(d^2)$
- Computation for implicit kernel calculation: $O(d)$.

Polynomial Kernel in \mathbb{R}^d

- Input space $\mathcal{X} = \mathbb{R}^d$
- Kernel function:

$$k(x, x') = \underline{(1 + \langle x, x' \rangle)}^M$$

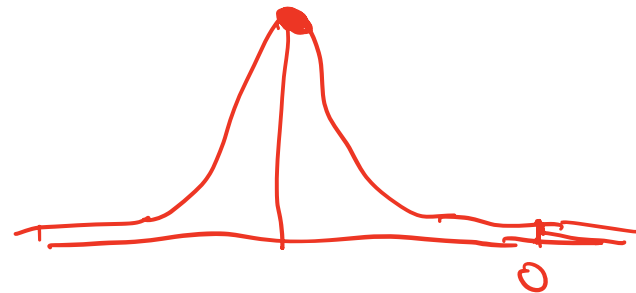
- Corresponds to a feature map with all monomials up to degree M .
- For any M , computing the kernel has same computational cost
- Cost of explicit inner product computation grows rapidly in M .

Radial Basis Function (RBF) / Gaussian Kernel

Input space $\mathcal{X} = \mathbb{R}^d$

$$k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

vector



where σ^2 is known as the bandwidth parameter.

- Probably the most common nonlinear kernel.

scalar.

$$k = k(x_i, x_j)$$

i, j

Radial Basis Function (RBF) / Gaussian Kernel

Input space $\mathcal{X} = \mathbb{R}^d$

$$k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right),$$

where σ^2 is known as the bandwidth parameter.

- Probably the most common nonlinear kernel.
- Does it act like a similarity score?

Radial Basis Function (RBF) / Gaussian Kernel

Input space $\mathcal{X} = \mathbb{R}^d$

$$k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right),$$

where σ^2 is known as the bandwidth parameter.

- Probably the most common nonlinear kernel.
- Does it act like a similarity score?
- Have we departed from our “inner product of feature vector” recipe?
 - Yes and no: corresponds to an infinite dimensional feature vector

Remaining Questions

Our current recipe:

- Recognize kernelized problem: $\psi(x)$ only occur in inner products $\psi(x)^T \psi(x')$

Remaining Questions

Our current recipe:

- Recognize kernelized problem: $\psi(x)$ only occur in inner products $\psi(x)^T \psi(x')$
- Pick a kernel function (“similarity score”)

Remaining Questions

Our current recipe:

- Recognize kernelized problem: $\psi(x)$ only occur in inner products $\psi(x)^T \psi(x')$
- Pick a kernel function (“similarity score”)
- Compute the kernel matrix (n by n where n is the dataset size)

Remaining Questions

Our current recipe:

- Recognize kernelized problem: $\psi(x)$ only occur in inner products $\psi(x)^T \psi(x')$
- Pick a kernel function (“similarity score”)
- Compute the kernel matrix (n by n where n is the dataset size)
- Optimize the model and make predictions by accessing the kernel matrix

Next: When can we apply kernelization?

SVM solution is in the “span of the data”

- We found the SVM dual problem can be written as:

$$\sup_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_j^T x_i$$

$$\text{s.t.} \quad \sum_{i=1}^n \alpha_i y_i = 0$$

$$\alpha_i \in \left[0, \frac{C}{n}\right] \quad i = 1, \dots, n.$$

- Given dual solution α^* , primal solution is $w^* = \sum_{i=1}^n \alpha_i^* y_i x_i$.

SVM solution is in the “span of the data”

- We found the SVM dual problem can be written as:

$$\begin{aligned} \sup_{\alpha \in \mathbb{R}^n} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_j^T x_i \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \\ & \alpha_i \in \left[0, \frac{C}{n}\right] \quad i = 1, \dots, n. \end{aligned}$$

- Given dual solution α^* , primal solution is $w^* = \sum_{i=1}^n \alpha_i^* y_i x_i$.
- Notice: w^* is a linear combination of training inputs x_1, \dots, x_n .

SVM solution is in the “span of the data”

- We found the SVM dual problem can be written as:

$$\begin{aligned} \sup_{\alpha \in \mathbb{R}^n} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_j^T x_i \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \\ & \alpha_i \in \left[0, \frac{C}{n}\right] \quad i = 1, \dots, n. \end{aligned}$$

- Given dual solution α^* , primal solution is $w^* = \sum_{i=1}^n \alpha_i^* y_i x_i$.
- Notice: w^* is a linear combination of training inputs x_1, \dots, x_n .
- We refer to this phenomenon by saying “ w^* is in the **span of the data.**”
 - Or in math, $w^* \in \text{span}(x_1, \dots, x_n)$.

Ridge regression solution is in the “span of the data”

- The ridge regression solution for regularization parameter $\lambda > 0$ is

$$w^* = \arg \min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \{w^T x_i - y_i\}^2 + \lambda \|w\|_2^2.$$

MSE

regul arization.

pred *target.*

Ridge regression solution is in the “span of the data”

- The ridge regression solution for regularization parameter $\lambda > 0$ is

$$w^* = \arg \min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \{w^T x_i - y_i\}^2 + \lambda \|w\|_2^2.$$

- This has a closed form solution:

$$w^* = (X^T X + \lambda I)^{-1} X^T y,$$

where X is the design matrix, with x_1, \dots, x_n as rows.

Ridge regression solution is in the “span of the data”

- Rearranging $w^* = (X^T X + \lambda I)^{-1} X^T y$, we can show that:

Ridge regression solution is in the “span of the data”

$$(A+B)^{-1}$$

- Rearranging $w^* = (X^T X + \lambda I)^{-1} X^T y$, we can show that:

$$\begin{aligned} w^* &= X^T \left(\frac{1}{\lambda} y - \frac{1}{\lambda} X w^* \right) \\ &= X^T \alpha^* = \sum_{i=1}^n \alpha_i^* x_i. \end{aligned}$$

- So w^* is in the span of the data
 - i.e. $w^* \in \text{span}(x_1, \dots, x_n)$

If solution is in the span of the data, we can reparameterize

- The ridge regression solution for regularization parameter $\lambda > 0$ is

$$w^* = \arg \min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \{w^T x_i - y_i\}^2 + \lambda \|w\|_2^2.$$

If solution is in the span of the data, we can reparameterize

- The ridge regression solution for regularization parameter $\lambda > 0$ is

$$w^* = \arg \min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \{w^T x_i - y_i\}^2 + \lambda \|w\|_2^2.$$

- We now know that $w^* \in \text{span}(x_1, \dots, x_n) \subset \mathbb{R}^d$.
- So rather than minimizing over all of \mathbb{R}^d , we can minimize over $\text{span}(x_1, \dots, x_n)$.

$$w^* = \arg \min_{w \in \text{span}(x_1, \dots, x_n)} \frac{1}{n} \sum_{i=1}^n \{w^T x_i - y_i\}^2 + \lambda \|w\|_2^2.$$

If solution is in the span of the data, we can reparameterize

- The ridge regression solution for regularization parameter $\lambda > 0$ is

$$w^* = \arg \min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \{w^T x_i - y_i\}^2 + \lambda \|w\|_2^2.$$

- We now know that $w^* \in \text{span}(x_1, \dots, x_n) \subset \mathbb{R}^d$.
- So rather than minimizing over all of \mathbb{R}^d , we can minimize over $\text{span}(x_1, \dots, x_n)$.

$$w^* = \arg \min_{w \in \text{span}(x_1, \dots, x_n)} \frac{1}{n} \sum_{i=1}^n \{w^T x_i - y_i\}^2 + \lambda \|w\|_2^2.$$

- Let's reparameterize the objective by replacing w as a linear combination of the inputs.

If solution is in the span of the data, we can reparameterize

- Note that for any $w \in \text{span}(x_1, \dots, x_n)$, we have $w = X^T \alpha$, for some $\alpha \in \mathbb{R}^n$.
- So let's replace w with $X^T \alpha$ in our optimization problem:

$$\text{[original]} \quad w^* = \arg \min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \{w^T x_i - y_i\}^2 + \lambda \|w\|_2^2$$

If solution is in the span of the data, we can reparameterize

- Note that for any $w \in \text{span}(x_1, \dots, x_n)$, we have $w = X^T \alpha$, for some $\alpha \in \mathbb{R}^n$.
- So let's replace w with $X^T \alpha$ in our optimization problem:

$$\text{[original]} \quad w^* = \arg \min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \{w^T x_i - y_i\}^2 + \lambda \|w\|_2^2$$

$$\text{[reparameterized]} \quad \alpha^* = \arg \min_{\alpha \in \mathbb{R}^n} \frac{1}{n} \sum_{i=1}^n \{ (X^T \alpha)^T x_i - y_i \}^2 + \lambda \|X^T \alpha\|_2^2.$$

- To get w^* from the reparameterized optimization problem, we just take $w^* = X^T \alpha^*$.
- We changed the dimension of our optimization variable from d to n . Is this useful?

Consider very large feature spaces

- Suppose we have a 300-million dimension feature space [very large]
 - (e.g. using high order monomial interaction terms as features, as described last lecture)
- Suppose we have a training set of 300,000 examples [fairly large]
- In the original formulation, we solve a 300-million dimension optimization problem.
- In the reparameterized formulation, we solve a 300,000-dimension optimization problem.
- This is why we care about when the solution is in the span of the data.
- This reparameterization is interesting when we have more features than data ($d \gg n$).

More General

- For SVM and ridge regression, we found that the solution is in the span of the data.
- The Representer Theorem shows that this “span of the data” result occurs far more generally.

The Representer Theorem (Optional)

- Generalized objective:

$$w^* = \arg \min_{w \in \mathcal{H}} R(\|w\|) + L(\langle w, x_1 \rangle, \dots, \langle w, x_n \rangle)$$

- Representer theorem tells us we can look for w^* in the span of the data:

$$w^* = \arg \min_{w \in \text{span}(x_1, \dots, x_n)} R(\|w\|) + L(\langle w, x_1 \rangle, \dots, \langle w, x_n \rangle).$$

- So we can reparameterize as before:

$$\alpha^* = \arg \min_{\alpha \in \mathbb{R}^n} R\left(\left\|\sum_{i=1}^n \alpha_i x_i\right\|\right) + L\left(\left\langle \sum_{i=1}^n \alpha_i x_i, x_1 \right\rangle, \dots, \left\langle \sum_{i=1}^n \alpha_i x_i, x_n \right\rangle\right).$$

- Our reparameterization trick applies much more broadly than SVM and ridge.

Summary

- We formulate the kernelized versions of SVM and ridge regression.

Summary

- We formulate the kernelized versions of SVM and ridge regression.
- Many other algorithms can be kernelized.

Summary

- We formulate the kernelized versions of SVM and ridge regression.
- Many other algorithms can be kernelized.
- Our principled tool for kernelization is reparameterization by the representer theorem.

Summary

- We formulate the kernelized versions of SVM and ridge regression.
- Many other algorithms can be kernelized.
- Our principled tool for kernelization is reparameterization by the representer theorem.
- Representer theorem says that all norm-regularized linear models can be kernelized.

Summary

- We formulate the kernelized versions of SVM and ridge regression.
- Many other algorithms can be kernelized.
- Our principled tool for kernelization is reparameterization by the representer theorem.
- Representer theorem says that all norm-regularized linear models can be kernelized.
- Once kernelized, we can apply the kernel trick: doesn't need to represent $\phi(x)$ explicitly.

Probabilistic Modeling: Overview

Why probabilistic modeling?

- A unified framework that covers many models, e.g., linear regression, logistic regression
- Learning as **statistical inference**
- Principled ways to incorporate your belief on the data generating distribution (inductive biases)

Two ways of generating data

- Two ways to model how the data is generated:

Two ways of generating data

- Two ways to model how the data is generated:
 - Conditional: $p(y | x)$ \rightarrow *classification*
 - Generative: $p(x, y)$

Two ways of generating data

- Two ways to model how the data is generated:
 - **Conditional:** $p(y | x)$
 - **Generative:** $p(x, y)$
- How to estimate the parameters of our model? Maximum likelihood estimation.

maximize likelihood. through finding the best parameters.

Two ways of generating data

- Two ways to model how the data is generated:
 - **Conditional:** $p(y | x)$
 - **Generative:** $p(x, y)$
- How to estimate the parameters of our model? Maximum likelihood estimation.
- Compare and contrast conditional and generative models.

Conditional models

Linear regression

Linear regression is one of the most important methods in machine learning and statistics.

Goal: Predict a real-valued target y (also called response) from a vector of features x (also called covariates).

Linear regression

Linear regression is one of the most important methods in machine learning and statistics.

Goal: Predict a real-valued **target** y (also called response) from a vector of **features** x (also called covariates).

Examples:

- Predicting house price given location, condition, build year etc.
- Predicting medical cost of a person given age, sex, region, BMI etc.
- Predicting age of a person based on their photos.

Problem setup

Data Training examples $\mathcal{D} = \{(x^{(n)}, y^{(n)})\}_{n=1}^N$, where $x \in \mathbb{R}^d$ and $y \in \mathbb{R}$.

Problem setup

Data Training examples $\mathcal{D} = \{(x^{(n)}, y^{(n)})\}_{n=1}^N$, where $x \in \mathbb{R}^d$ and $y \in \mathbb{R}$.

Model A *linear* function h (parametrized by θ) to predict y from x :

$$h(x) = \sum_{i=0}^d \theta_i x_i = \theta^T x, \quad (1)$$

(Handwritten red annotations: a red underline under $h(x)$, a red arrow pointing from θ^T to w^T , and a red 'x' over x)

where $\theta \in \mathbb{R}^d$ are the **parameters** (also called weights).

Problem setup

Data Training examples $\mathcal{D} = \{(x^{(n)}, y^{(n)})\}_{n=1}^N$, where $x \in \mathbb{R}^d$ and $y \in \mathbb{R}$.

Model A *linear* function h (parametrized by θ) to predict y from x :

$$h(x) = \sum_{i=0}^d \theta_i x_i = \theta^T x, \quad (1)$$

where $\theta \in \mathbb{R}^d$ are the **parameters** (also called weights).

$x^{(n)}$
 d
 j

Note that

- We incorporate the **bias term** (also called the intercept term) into x (i.e. $x_0 = 1$).
- We use superscript to denote the example id and subscript to denote the dimension id.

Parameter estimation

Loss function We estimate θ by minimizing the **squared loss** (the least square method):

$$\underline{\underline{J(\theta)}} = \frac{1}{N} \sum_{n=1}^N \underbrace{\left(y^{(n)} - \theta^T x^{(n)} \right)^2}_{\text{empirical risk}}. \quad (2)$$

Parameter estimation

Loss function We estimate θ by minimizing the **squared loss** (the least square method):

$$J(\theta) = \frac{1}{N} \sum_{n=1}^N \left(y^{(n)} - \theta^T x^{(n)} \right)^2. \quad (\text{empirical risk}) \quad (2)$$

Matrix form

- Let $X \in \mathbb{R}^{N \times d}$ be the design matrix whose rows are input features.
- Let $y \in \mathbb{R}^N$ be the vector of all targets.

Parameter estimation

Loss function We estimate θ by minimizing the **squared loss** (the least square method):

$$J(\theta) = \frac{1}{N} \sum_{n=1}^N \left(y^{(n)} - \theta^T x^{(n)} \right)^2. \quad (\text{empirical risk}) \quad (2)$$

Matrix form

- Let $X \in \mathbb{R}^{N \times d}$ be the **design matrix** whose rows are input features.
- Let $y \in \mathbb{R}^N$ be the vector of all targets.
- We want to solve

$$\hat{\theta} = \arg \min_{\theta} (X\theta - y)^T (X\theta - y). \quad (3)$$

Parameter estimation

Loss function We estimate θ by minimizing the **squared loss** (the least square method):

$$J(\theta) = \frac{1}{N} \sum_{n=1}^N \left(y^{(n)} - \theta^T x^{(n)} \right)^2. \quad (\text{empirical risk}) \quad (2)$$

- Matrix form**
- Let $X \in \mathbb{R}^{N \times d}$ be the **design matrix** whose rows are input features.
 - Let $y \in \mathbb{R}^N$ be the vector of all targets.
 - We want to solve

$$\hat{\theta} = \arg \min_{\theta} (X\theta - y)^T (X\theta - y). \quad (3)$$

Solution Closed-form solution: $\hat{\theta} = \underline{(X^T X)^{-1} X^T y}$.

Parameter estimation

Loss function We estimate θ by minimizing the **squared loss** (the least square method):

$$J(\theta) = \frac{1}{N} \sum_{n=1}^N \left(y^{(n)} - \theta^T x^{(n)} \right)^2. \quad \text{(empirical risk)} \quad (2)$$

→ Square.

- Matrix form**
- Let $X \in \mathbb{R}^{N \times d}$ be the **design matrix** whose rows are input features.
 - Let $y \in \mathbb{R}^N$ be the vector of all targets.
 - We want to solve

$$\hat{\theta} = \arg \min_{\theta} (X\theta - y)^T (X\theta - y). \quad (3)$$

Solution Closed-form solution: $\hat{\theta} = (X^T X)^{-1} X^T y$.

Review questions

- Derive the solution for linear regression.
- What if $X^T X$ is not invertible?

Review

We've seen

- Linear regression: response is a linear function of the inputs
- Estimate parameters by minimize the squared loss

Review

We've seen

- Linear regression: response is a linear function of the inputs
- Estimate parameters by minimize the squared loss

But...

- Why squared loss is a reasonable choice for regression problems?
- What assumptions are we making on the data? ([inductive bias](#))

Review

We've seen

- Linear regression: response is a linear function of the inputs
- Estimate parameters by minimize the squared loss

But...

- Why squared loss is a reasonable choice for regression problems?
- What assumptions are we making on the data? ([inductive bias](#))

Next,

- Derive linear regression from a [probabilistic modeling perspective](#).

Assumptions in linear regression

Assumptions in linear regression

- x and y are related through a linear function:

$$y = \theta^T x + \epsilon, \tag{4}$$

where ϵ is the **residual error** capturing all unmodeled effects (e.g., noise).

Assumptions in linear regression

- x and y are related through a linear function:

$$y = \theta^T x + \epsilon, \quad \leftarrow \text{residual error.} \quad (4)$$

where ϵ is the **residual error** capturing all unmodeled effects (e.g., noise).

- The errors are distributed *iid* (independently and identically distributed):

$$\epsilon \sim \mathcal{N}(0, \sigma^2). \quad (5)$$

Mean Variance.

Assumptions in linear regression

- x and y are related through a linear function:

$$y = \underline{\theta^T x} + \epsilon, \quad (4)$$

where ϵ is the **residual error** capturing all unmodeled effects (e.g., noise).

- The errors are distributed *iid* (independently and identically distributed):

$$\epsilon \sim \mathcal{N}(0, \sigma^2). \quad (5)$$

What's the distribution of $Y | X = x$?

$$Y \sim \mathcal{N}(\theta^T x, \sigma^2).$$

Assumptions in linear regression

- x and y are related through a linear function:

$$y = \theta^T x + \epsilon, \tag{4}$$

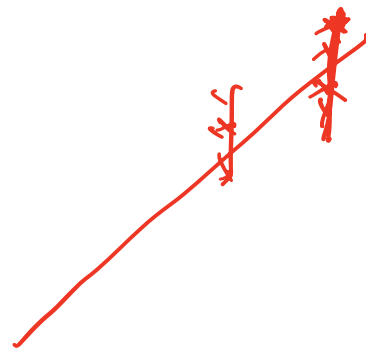
where ϵ is the **residual error** capturing all unmodeled effects (e.g., noise).

- The errors are distributed *iid* (independently and identically distributed):

$$\epsilon \sim \mathcal{N}(0, \sigma^2). \tag{5}$$

What's the distribution of $Y | X = x$?

$$p(y | x; \theta) = \mathcal{N}(\theta^T x, \sigma^2). \tag{6}$$



Imagine putting a Gaussian bump around the output of the linear predictor.

Assumptions in linear regression

- x and y are related through a linear function:

$$y = \theta^T x + \epsilon, \quad (4)$$

where ϵ is the **residual error** capturing all unmodeled effects (e.g., noise).

- The errors are distributed *iid* (independently and identically distributed):

$$\epsilon \sim \mathcal{N}(0, \sigma^2). \quad (5)$$

What's the distribution of $Y \mid X = x$?

$$p(y \mid x; \theta) = \mathcal{N}(\theta^T x, \sigma^2). \quad (6)$$

Imagine putting a Gaussian bump around the output of the linear predictor.

Maximum likelihood estimation (MLE)

Given a probabilistic model and a dataset \mathcal{D} , how to estimate the model parameters θ ?

(8)

Maximum likelihood estimation (MLE)

Given a probabilistic model and a dataset \mathcal{D} , how to estimate the model parameters θ ?

The **maximum likelihood principle** says that we should maximize the (conditional) likelihood of the data:

$$L(\theta) \stackrel{\text{def}}{=} p(\mathcal{D}; \theta) \tag{7}$$

(8)

Maximum likelihood estimation (MLE)

Given a probabilistic model and a dataset \mathcal{D} , how to estimate the model parameters θ ?

The **maximum likelihood principle** says that we should maximize the (conditional) likelihood of the data:

likelihood is a function of θ .

$$L(\theta) \stackrel{\text{def}}{=} p(\mathcal{D}; \theta) \quad \textcircled{x} \ y. \quad y|x_j \quad (7)$$

$$= \prod_{n=1}^N p(y^{(n)} | x^{(n)}; \theta). \quad \text{(examples are distributed } iid) \quad (8)$$

Maximum likelihood estimation (MLE)

Given a probabilistic model and a dataset \mathcal{D} , how to estimate the model parameters θ ?

The **maximum likelihood principle** says that we should maximize the (conditional) likelihood of the data:

$$L(\theta) \stackrel{\text{def}}{=} p(\mathcal{D}; \theta) \quad (7)$$

$$\log \left(\prod_{n=1}^N p(y^{(n)} | x^{(n)}; \theta) \right) \quad (\text{examples are distributed } iid) \quad (8)$$

$\rightarrow \sum \log p(y^{(i)} | x^{(i)})$

In practice, we maximize the log likelihood $\ell(\theta)$, or equivalently, minimize the negative log likelihood (NLL).

MLE for linear regression

Let's find the MLE solution for our model. Recall that $Y | X = x \sim \mathcal{N}(\theta^T x, \sigma^2)$.

$$\begin{aligned} \text{log likelihood} \leftarrow L(\theta) &= \log L(\theta) \\ &= \log \prod_n p(y^{(n)} | x^{(n)}; \theta) \\ &= \sum_n \log \underbrace{p(y^{(n)} | x^{(n)}; \theta)}_{\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(n)} - \theta^T x^{(n)})^2}{2\sigma^2}\right)} \cdot \mathcal{N}(\theta^T x, \sigma^2). \end{aligned}$$

max.

$$= N \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} \sum_n (y^{(n)} - \theta^T x^{(n)})^2$$

divide by $N \rightarrow C - \underbrace{\frac{1}{2\sigma^2} \frac{1}{N} \sum_n (y^{(n)} - \theta^T x^{(n)})^2}_{\text{MSE}}$

Gradient of the likelihood

Recall that we obtained the normal equation by setting the derivative of the squared loss to zero. Now let's compute the derivative of the likelihood w.r.t. the parameters.

$$L(\theta) = N \log \frac{1}{\sqrt{2\pi} \sigma} - \frac{1}{2\sigma^2} \sum_n (y^{(n)} - \theta^T x^{(n)})^2$$
$$\frac{\partial L(\theta)}{\partial \theta_i} = -\frac{1}{\sigma^2} \sum_n (y^{(n)} - \theta^T x^{(n)}) \cdot x_i^{(n)}$$
$$= -\frac{1}{\sigma^2} \sum_n \left(\underset{\substack{\downarrow \\ \text{target}}}{y^{(n)}} - \underset{\substack{\downarrow \\ \text{pred}}}{\hat{y}^{(n)}} \right) \cdot \underset{\substack{\downarrow \\ \text{input}}}{x_i^{(n)}}$$

Review

We've seen

- Linear regression assumes that $Y|X=x$ follows a Gaussian distribution
- MLE of linear regression is equivalent to the least square method

Review

We've seen

- Linear regression assumes that $Y | X = x$ follows a Gaussian distribution
- MLE of linear regression is equivalent to the least square method

However,

$$Y = \begin{pmatrix} +1 / -1 \\ 0 / 1 \end{pmatrix}$$

- Sometimes Gaussian distribution is not a reasonable assumption, e.g., classification
- Can we use the same modeling approach for other prediction tasks?

Review

We've seen

- Linear regression assumes that $Y | X = x$ follows a Gaussian distribution
- MLE of linear regression is equivalent to the least square method

However,

- Sometimes Gaussian distribution is not a reasonable assumption, e.g., classification
- Can we use the same modeling approach for other prediction tasks?

Next,

- Derive [logistic regression](#) for classification.

Assumptions in logistic regression

Consider binary classification where $Y \in \{0, 1\}$. What should be the distribution $Y | X = x$?

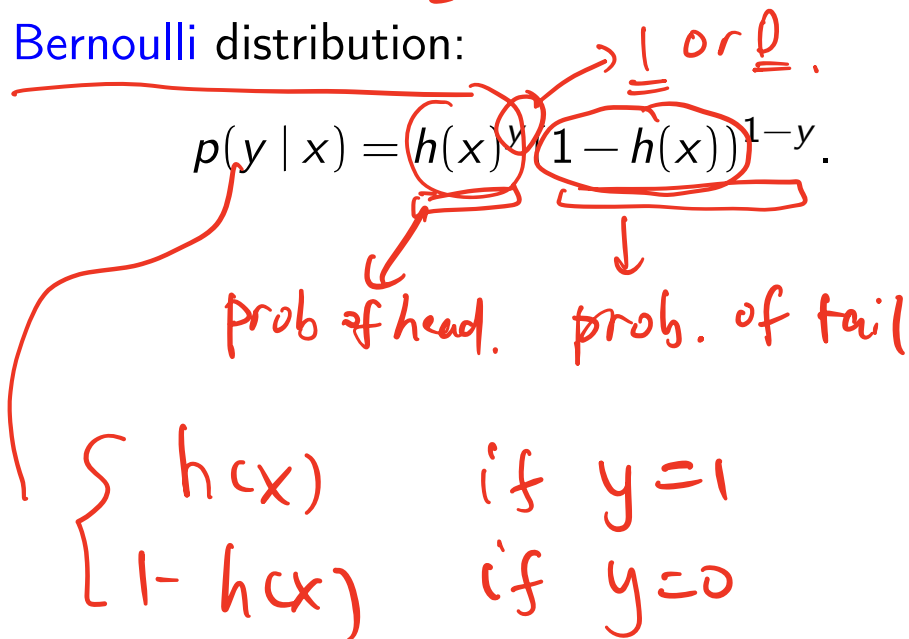
Assumptions in logistic regression

Consider binary classification where $Y \in \{0, 1\}$. What should be the distribution $Y | X = x$?

We model $p(y | x)$ as a **Bernoulli** distribution:

$$p(y | x) = h(x)^y (1 - h(x))^{1-y}.$$

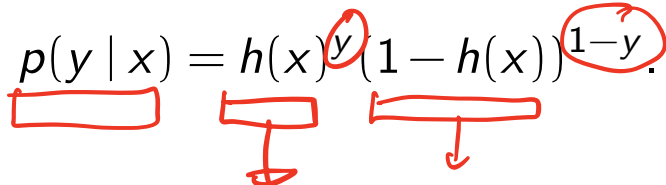
(9)



Assumptions in logistic regression

Consider binary classification where $Y \in \{0, 1\}$. What should be the distribution $Y | X = x$?

We model $p(y | x)$ as a **Bernoulli** distribution:

$$p(y | x) = h(x)^y (1 - h(x))^{1-y} \quad (9)$$


How should we parameterize $h(x)$?

↳ prob. of $y = 1$.


Assumptions in logistic regression

Consider binary classification where $Y \in \{0, 1\}$. What should be the distribution $Y | X = x$?

We model $p(y | x)$ as a **Bernoulli** distribution:

$$p(y | x) = h(x)^y (1 - h(x))^{1-y}. \quad (9)$$

How should we parameterize $h(x)$?

- What is $p(y = 1 | x)$ and $p(y = 0 | x)$?


Assumptions in logistic regression

Consider binary classification where $Y \in \{0, 1\}$. What should be the distribution $Y | X = x$?

We model $p(y | x)$ as a **Bernoulli** distribution:

$$p(y | x) = h(x)^y (1 - h(x))^{1-y}. \quad (9)$$

How should we parameterize $h(x)$?

- What is $p(y = 1 | x)$ and $p(y = 0 | x)$? $h(x) \in (0, 1)$.

Assumptions in logistic regression

Consider binary classification where $Y \in \{0, 1\}$. What should be the distribution $Y | X = x$?

We model $p(y | x)$ as a **Bernoulli** distribution:

$$p(y | x) = h(x)^y (1 - h(x))^{1-y}. \quad (9)$$

How should we parameterize $h(x)$?

- What is $p(y = 1 | x)$ and $p(y = 0 | x)$? $h(x) \in (0, 1)$.
- What is the mean of $Y | X = x$?

Assumptions in logistic regression

Consider binary classification where $Y \in \{0, 1\}$. What should be the distribution $Y | X = x$?

We model $p(y | x)$ as a **Bernoulli** distribution:

$$p(y | x) = h(x)^y (1 - h(x))^{1-y}. \quad (9)$$

How should we parameterize $h(x)$?

- What is $p(y = 1 | x)$ and $p(y = 0 | x)$? $h(x) \in (0, 1)$.
- What is the mean of $Y | X = x$? $h(x)$. (Think how we parameterize the mean in linear regression)

Assumptions in logistic regression

Consider binary classification where $Y \in \{0, 1\}$. What should be the distribution $Y | X = x$?

We model $p(y | x)$ as a **Bernoulli** distribution:

$$p(y | x) = h(x)^y (1 - h(x))^{1-y}. \quad (9)$$

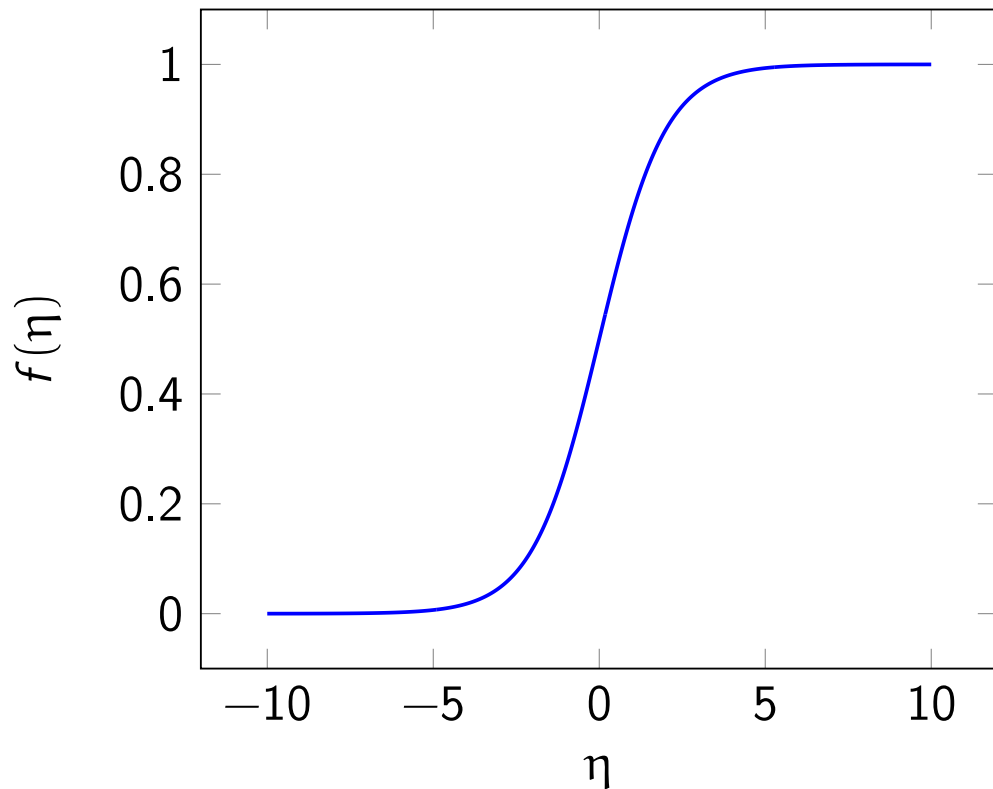
How should we parameterize $h(x)$?

- What is $p(y = 1 | x)$ and $p(y = 0 | x)$? $h(x) \in (0, 1)$.
- What is the mean of $Y | X = x$? $h(x)$. (Think how we parameterize the mean in linear regression)
- Need a function f to map the linear predictor $\theta^T x$ in \mathbb{R} to $(0, 1)$:

$$f(\eta) = \frac{1}{1 + e^{-\eta}} \quad \text{logistic function} \quad (10)$$

Logistic regression

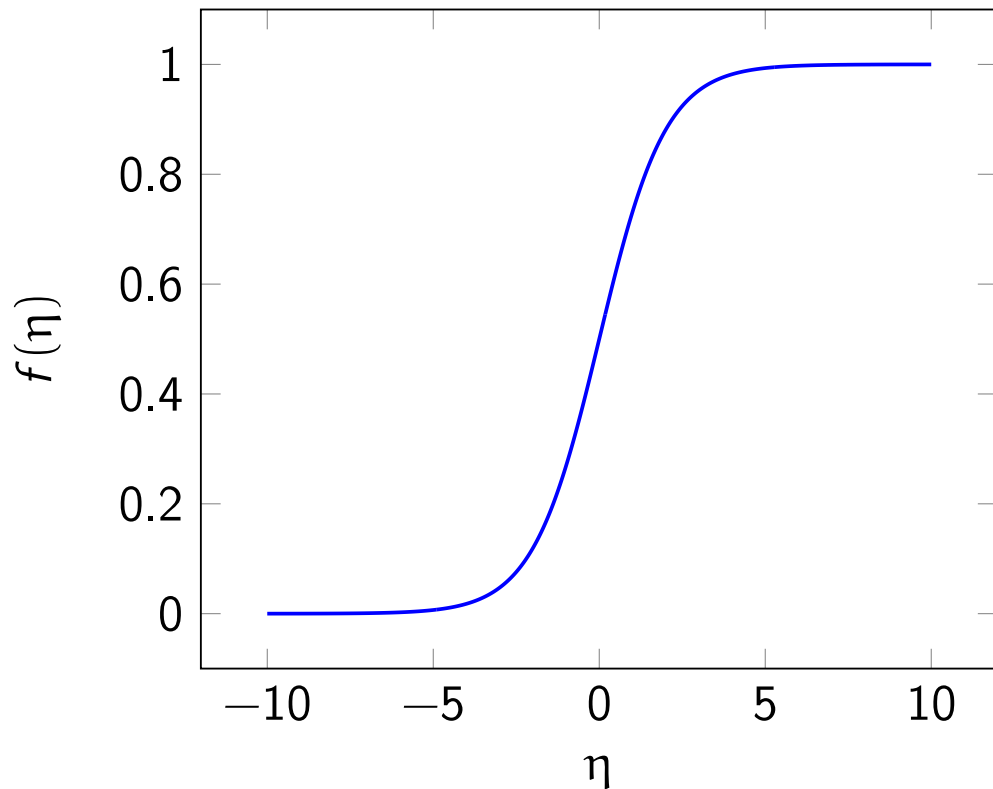
$$f(\eta) = \frac{1}{1+e^{-\eta}}$$



- $p(y | x) = \text{Bernoulli}(f(\theta^T x))$.

Logistic regression

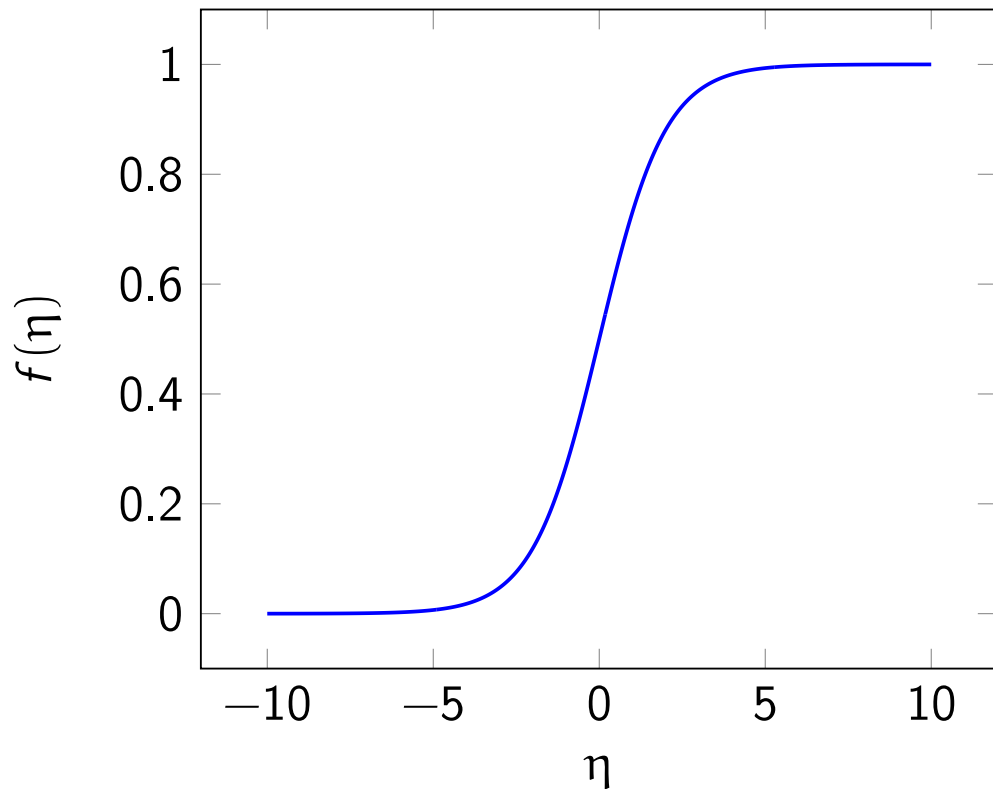
$$f(\eta) = \frac{1}{1+e^{-\eta}}$$



- $p(y | x) = \text{Bernoulli}(f(\theta^T x))$.
- When do we have $p(y = 1 | x) = 1$ and $p(y = 0 | x) = 1$?

Logistic regression

$$f(\eta) = \frac{1}{1+e^{-\eta}}$$

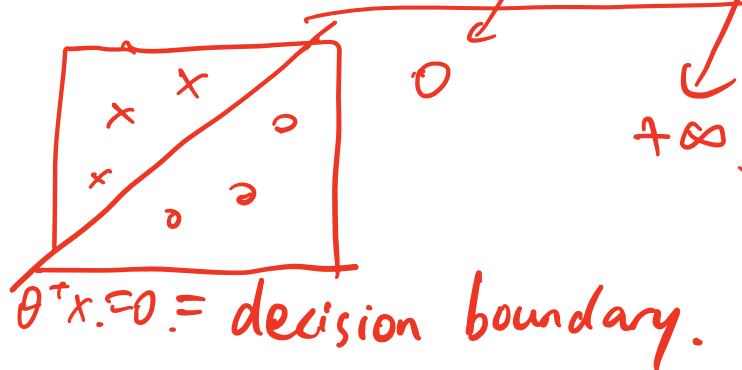


- $p(y | x) = \text{Bernoulli}(f(\theta^T x))$.
- When do we have $p(y = 1 | x) = 1$ and $p(y = 0 | x) = 1$?
- **Exercise:** show that the **log odds** is

$$\log \frac{p(y = 1 | x)}{p(y = 0 | x)} = \theta^T x. \quad (11)$$

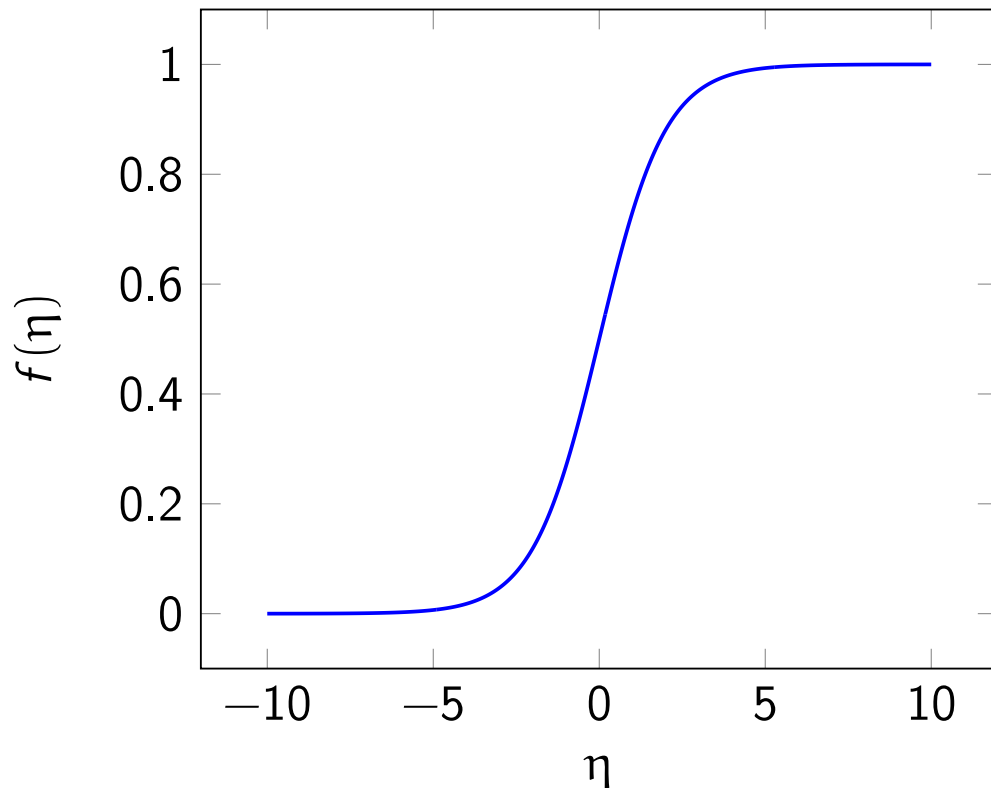
Handwritten notes: "big" with an arrow pointing to the fraction, and a circle around $\theta^T x$.

$$\Rightarrow \text{linear decision boundary} \quad (12)$$



Logistic regression

$$f(\eta) = \frac{1}{1+e^{-\eta}}$$



- $p(y | x) = \text{Bernoulli}(f(\theta^T x))$.
- When do we have $p(y = 1 | x) = 1$ and $p(y = 0 | x) = 1$?
- **Exercise:** show that the **log odds** is

$$\log \frac{p(y = 1 | x)}{p(y = 0 | x)} = \theta^T x. \quad (11)$$

$$\implies \text{linear decision boundary} \quad (12)$$

- How do we extend it to multiclass classification? (more on this later)