

Kernels & Probabilistic Modeling

Mengye Ren

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

NYU

October 1, 2023

- Oct 15: Legislative Day No Class
- Oct 15: Homework 2 Due
- Oct 22: Midterm, in class, covers everything including next week's content (Oct 8)

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):

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

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)

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

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

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

- If α^* is an optimal value, then

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

SVM Dual Problem

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

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

- If α^* is an optimal value, then

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

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

Compute the Inner Products

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

$$\begin{aligned}\psi(x)^T \psi(x') &= x_1^2 x_1'^2 + (\sqrt{2}x_1x_2)(\sqrt{2}x_1'x_2') + x_2^2 x_2'^2 \\ &= (x_1x_1')^2 + 2(x_1x_1')(x_2x_2') + (x_2x_2')^2 \\ &= (x_1x_1' + x_2x_2')^2 \\ &= (x^T x')^2\end{aligned}$$

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

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

- Using explicit features: $O(d^p)$
- Using implicit computation: $O(d)$

Kernel Function

The Kernel Function

- **Input space:** \mathcal{X}
- **Feature space:** \mathcal{H} (a Hilbert space, e.g. \mathbb{R}^d)
- **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} (a Hilbert space, e.g. \mathbb{R}^d)
- **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 $k(x, x')$?

The Kernel Function

- **Input space:** \mathcal{X}
- **Feature space:** \mathcal{H} (a Hilbert space, e.g. \mathbb{R}^d)
- **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} (a Hilbert space, e.g. \mathbb{R}^d)
- **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.

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

Definition

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

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

- 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} :

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

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.

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 $d \gg n$.

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

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

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:

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

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:

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

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

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:

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

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

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

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

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:

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

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

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

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

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

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

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

$$\psi(x) = x$$

- Kernel:

$$k(x, x') = x^T x'$$

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:

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

- 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:

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

- 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 $\mathcal{X} = \mathbb{R}^d$

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

where σ^2 is known as the bandwidth parameter.

- Probably the most common nonlinear kernel.

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:

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

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

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

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

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:

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

- 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 ($d \gg n$).

- 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.

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

Summary

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

Summary

- We formalize 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 formalized 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 formalized 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:** $p(y | x)$
 - **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.

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 :

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

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).

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):

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

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

Parameter estimation

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

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

- 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?

We've seen

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

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**)

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:

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

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

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

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

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 $Y \mid X = x$?

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.

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:

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

$$= \prod_{n=1}^N p(y^{(n)} \mid x^{(n)}; \theta). \tag{8}$$

(examples are distributed *iid*)

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:

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

$$= \prod_{n=1}^N p(y^{(n)} \mid x^{(n)}; \theta). \tag{8}$$

(examples are distributed *iid*)

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

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.

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

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?

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:

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

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 $h(x)$?

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:

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

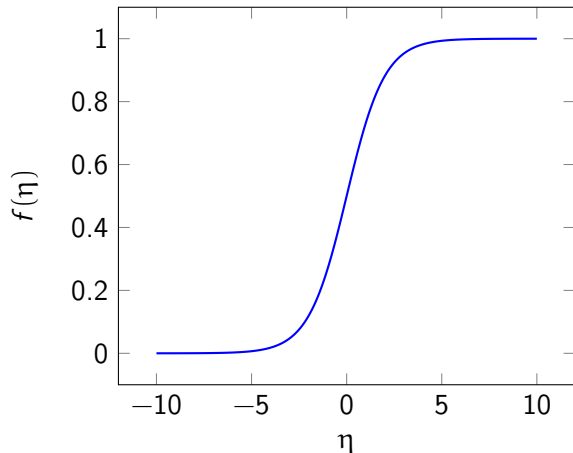
How should we parameterize $h(x)$?

- 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 $(0, 1)$:

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

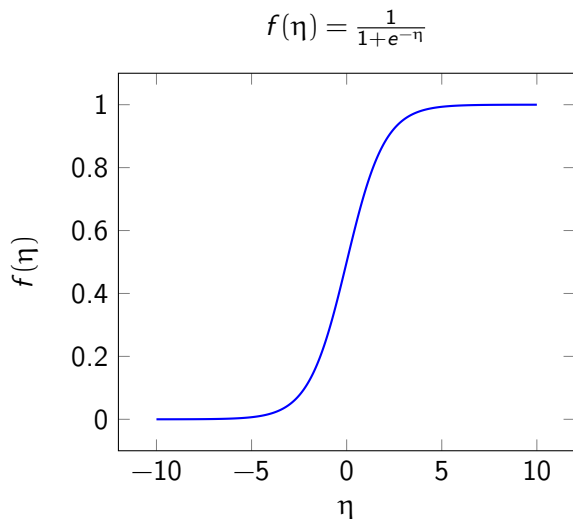
Logistic regression

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



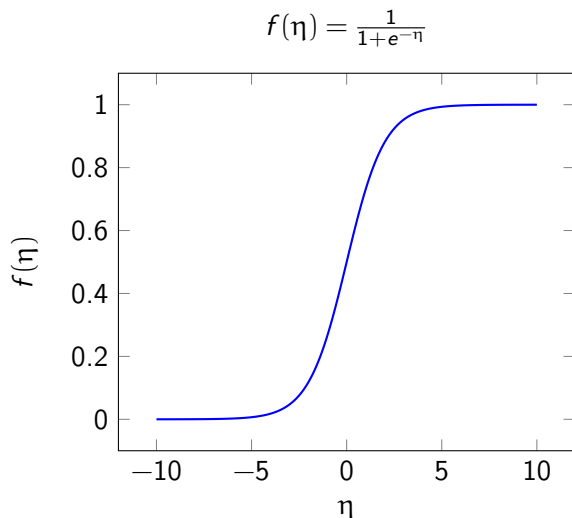
- $p(y | x) = \text{Bernoulli}(f(\theta^T x)).$

Logistic regression



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

Logistic regression

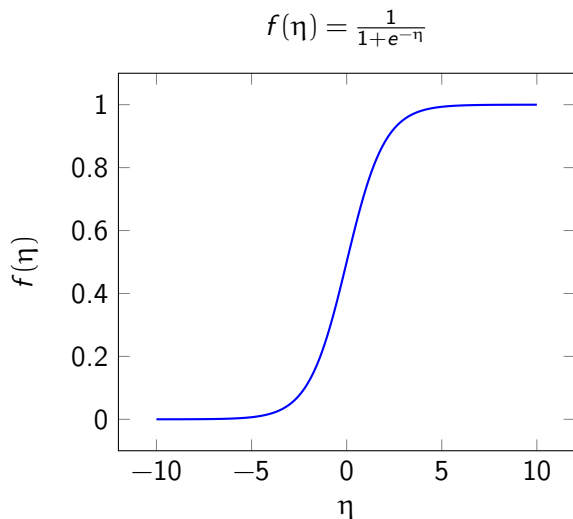


- $p(y | x) = \text{Bernoulli}(f(\theta^T x))$.
- When do we have $p(y = 1 | x) = 1$ and $p(y = 0 | x) = 1$?
- **Exercise:** show that the **log odds** is

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

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

Logistic regression



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

MLE for logistic regression

Similar to linear regression, let's estimate θ by maximizing the conditional log likelihood.

MLE for logistic regression

Similar to linear regression, let's estimate θ by maximizing the conditional log likelihood.

- Closed-form solutions are not available.
- But, the likelihood is concave—**gradient ascent** gives us the unique optimal solution.

$$\theta := \theta + \alpha \nabla_{\theta} \ell(\theta). \quad (13)$$

Gradient descent for logistic regression

Math review: Chain rule

If z depends on y which itself depends on x , e.g., $z = (y(x))^2$, then $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$.

Gradient descent for logistic regression

Math review: Chain rule

If z depends on y which itself depends on x , e.g., $z = (y(x))^2$, then $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$.

Likelihood for a single example: $\ell^n = y^{(n)} \log f(\theta^T x^{(n)}) + (1 - y^{(n)}) \log(1 - f(\theta^T x^{(n)}))$.

A closer look at the gradient

$$\frac{\partial \ell}{\partial \theta_i} = \sum_{n=1}^N (y^{(n)} - f(\theta^T x^{(n)})) x_i^{(n)} \quad (14)$$

A closer look at the gradient

$$\frac{\partial \ell}{\partial \theta_i} = \sum_{n=1}^N (y^{(n)} - f(\theta^T x^{(n)})) x_i^{(n)} \quad (14)$$

- Does this look familiar?
- Our derivation for linear regression and logistic regression are quite similar...
- Next, a more general family of models.

Compare linear regression and logistic regression

linear regression	logistic regression
-------------------	---------------------

Compare linear regression and logistic regression

	linear regression	logistic regression
Combine the inputs	$\theta^T x$ (linear)	$\theta^T x$ (linear)

Compare linear regression and logistic regression

	linear regression	logistic regression
Combine the inputs	$\theta^T x$ (linear)	$\theta^T x$ (linear)
Output	real	categorical
Conditional distribution	Gaussian	Bernoulli

Compare linear regression and logistic regression

	linear regression	logistic regression
Combine the inputs	$\theta^T x$ (linear)	$\theta^T x$ (linear)
Output	real	categorical
Conditional distribution	Gaussian	Bernoulli
Transfer function $f(\theta^T x)$	identity	logistic

Compare linear regression and logistic regression

	linear regression	logistic regression
Combine the inputs	$\theta^T x$ (linear)	$\theta^T x$ (linear)
Output	real	categorical
Conditional distribution	Gaussian	Bernoulli
Transfer function $f(\theta^T x)$	identity	logistic
Mean $\mathbb{E}(Y X = x; \theta)$	$f(\theta^T x)$	$f(\theta^T x)$

Compare linear regression and logistic regression

	linear regression	logistic regression
Combine the inputs	$\theta^T x$ (linear)	$\theta^T x$ (linear)
Output	real	categorical
Conditional distribution	Gaussian	Bernoulli
Transfer function $f(\theta^T x)$	identity	logistic
Mean $\mathbb{E}(Y X = x; \theta)$	$f(\theta^T x)$	$f(\theta^T x)$

- x enters through a linear function.
- The main difference between the formulations is due to different conditional distributions.
- Can we generalize the idea to handle other output types, e.g., positive integers?

Construct a generalized regression model

Task: Given x , predict $p(y | x)$

Modeling:

- Choose a parametric family of distributions $p(y; \theta)$ with parameters $\theta \in \Theta$
- Choose a transfer function that maps a linear predictor in \mathbb{R} to Θ

$$\underbrace{x}_{\in \mathbb{R}^d} \mapsto \underbrace{w^T x}_{\in \mathbb{R}} \mapsto \underbrace{f(w^T x)}_{\in \Theta} = \theta, \quad (15)$$

Learning: MLE: $\hat{\theta} \in \arg \max_{\theta} \log p(\mathcal{D}; \theta)$

Inference: For prediction, use $x \rightarrow f(w^T x)$

Construct a generalized regression model

Task: Given x , predict $p(y | x)$

Modeling:

- Choose a parametric family of distributions $p(y; \theta)$ with parameters $\theta \in \Theta$
- Choose a transfer function that maps a linear predictor in \mathbb{R} to Θ

$$\underbrace{x}_{\in \mathbb{R}^d} \mapsto \underbrace{w^T x}_{\in \mathbb{R}} \mapsto \underbrace{f(w^T x)}_{\in \Theta} = \theta, \quad (15)$$

Learning: MLE: $\hat{\theta} \in \arg \max_{\theta} \log p(\mathcal{D}; \theta)$

Inference: For prediction, use $x \rightarrow f(w^T x)$

Example: Construct Poisson regression

Say we want to predict the number of people entering a restaurant in New York during lunch time.

- What features would be useful?
- What's a good model for number of visitors (the **output distribution**)?

Example: Construct Poisson regression

Say we want to predict the number of people entering a restaurant in New York during lunch time.

- What features would be useful?
- What's a good model for number of visitors (the **output distribution**)?

Math review: Poisson distribution

Given a random variable $Y \in 0, 1, 2, \dots$ following $\text{Poisson}(\lambda)$, we have

$$p(Y = k; \lambda) = \frac{\lambda^k e^{-\lambda}}{k!}, \quad (16)$$

where $\lambda > 0$ and $\mathbb{E}[Y] = \lambda$.

The Poisson distribution is usually used to model the number of events occurring during a fixed period of time.

Example: Construct Poisson regression

We've decided that $Y | X = x \sim \text{Poisson}(\eta)$, what should be the transfer function f ?
 x enters linearly:

$$x \mapsto \underbrace{w^T x}_R \mapsto \lambda = \underbrace{f(w^T x)}_{(0, \infty)}$$

Standard approach is to take

$$f(w^T x) = \exp(w^T x).$$

Likelihood of the full dataset $\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}$:

$$\log p(y_i; \lambda_i) = [y_i \log \lambda_i - \lambda_i - \log(y_i!)] \quad (17)$$

$$\log p(\mathcal{D}; w) = \sum_{i=1}^n [y_i \log [\exp(w^T x_i)] - \exp(w^T x_i) - \log(y_i!)] \quad (18)$$

$$= \sum_{i=1}^n [y_i w^T x_i - \exp(w^T x_i) - \log(y_i!)] \quad (19)$$

Example: multinomial logistic regression

How to extend logistic regression to multiclass classification?

Example: multinomial logistic regression

How to extend logistic regression to multiclass classification? Output: Bernoulli distribution \rightarrow **categorical distribution**

- Parametrized by a probability vector $\theta = (\theta_1, \dots, \theta_k) \in \mathbb{R}^k$:
 - $\sum_{i=1}^k \theta_i = 1$ and $\theta_i \geq 0$ for $i = 1, \dots, k$
 - So $\forall y \in \{1, \dots, k\}$, $p(y) = \theta_y$.
- From each x , we compute a linear score function for each class:

$$x \mapsto (\langle w_1, x \rangle, \dots, \langle w_k, x \rangle) \in \mathbb{R}^k,$$

- What's the transfer function that maps this \mathbb{R}^k vector into a probability?
The **softmax function**:

$$(s_1, \dots, s_k) \mapsto \theta = \left(\frac{e^{s_1}}{\sum_{i=1}^k e^{s_i}}, \dots, \frac{e^{s_k}}{\sum_{i=1}^k e^{s_i}} \right).$$

Recipe for constructing a conditional distribution for prediction:

- 1 Define input and output space (as for any other model).
- 2 Choose the output distribution $p(y | x; \theta)$ based on the task
- 3 Choose the transfer function that maps $w^T x$ to a Θ .
- 4 (The formal family is called “generalized linear models”.)

Learning:

- Fit the model by maximum likelihood estimation.
- Closed solutions do not exist in general, so we use gradient ascent.

Generative models

We've seen

- Model the conditional distribution $p(y | x; \theta)$ using generalized linear models.
- (Previously) Directly map x to y , e.g., perceptron.

Next,

- Model the **joint distribution** $p(x, y; \theta)$.
- Predict the label for x as $\arg \max_{y \in \mathcal{Y}} p(x, y; \theta)$.

Generative modeling through the Bayes rule

Training:

$$p(x, y) \tag{20}$$

(22)

Generative modeling through the Bayes rule

Training:

$$p(x, y) = p(x | y)p(y) \quad (20)$$

(22)

Generative modeling through the Bayes rule

Training:

$$p(x, y) = p(x | y)p(y) \quad (20)$$

(22)

Generative modeling through the Bayes rule

Training:

$$p(x, y) = p(x | y)p(y) \quad (20)$$

Testing:

$$p(y | x) \quad (22)$$

Generative modeling through the Bayes rule

Training:

$$p(x, y) = p(x | y)p(y) \quad (20)$$

Testing:

$$p(y | x) = \frac{p(x | y)p(y)}{p(x)} \quad \text{Bayes rule} \quad (21)$$

(22)

Generative modeling through the Bayes rule

Training:

$$p(x, y) = p(x | y)p(y) \quad (20)$$

Testing:

$$p(y | x) = \frac{p(x | y)p(y)}{p(x)} \quad \text{Bayes rule} \quad (21)$$

$$\arg \max_y p(y | x) = \arg \max_y p(x | y)p(y) \quad (22)$$

Naive Bayes (NB) models

Let's consider binary text classification (e.g., fake vs genuine review) as a motivating example.

Naive Bayes (NB) models

Let's consider binary text classification (e.g., fake vs genuine review) as a motivating example.

Bag-of-words representation of a document

- ["machine", "learning", "is", "fun", "."]
- $x_i \in \{0, 1\}$: whether the i -th word in our vocabulary exists in the input

$$x = [x_1, x_2, \dots, x_d] \quad \text{where } d = \text{vocabulary size} \quad (23)$$

Naive Bayes (NB) models

Let's consider binary text classification (e.g., fake vs genuine review) as a motivating example.

Bag-of-words representation of a document

- ["machine", "learning", "is", "fun", "."]
- $x_i \in \{0, 1\}$: whether the i -th word in our vocabulary exists in the input

$$x = [x_1, x_2, \dots, x_d] \quad \text{where } d = \text{vocabulary size} \quad (23)$$

What's the probability of a document x ?

Naive Bayes (NB) models

Let's consider binary text classification (e.g., fake vs genuine review) as a motivating example.

Bag-of-words representation of a document

- ["machine", "learning", "is", "fun", "."]
- $x_i \in \{0, 1\}$: whether the i -th word in our vocabulary exists in the input

$$x = [x_1, x_2, \dots, x_d] \quad \text{where } d = \text{vocabulary size} \quad (23)$$

What's the probability of a document x ?

$$p(x | y) = p(x_1, \dots, x_d | y) \quad (24)$$

$$= p(x_1 | y) p(x_2 | y, x_1) p(x_3 | y, x_2, x_1) \dots p(x_d | y, x_{d-1}, \dots, x_1) \quad \text{chain rule} \quad (25)$$

$$= \prod_{i=1}^d p(x_i | y, x_{<i}) \quad (26)$$

Naive Bayes assumption

Challenge: $p(x_i | y, x_{<i})$ is hard to model (and estimate), especially for large i .

Naive Bayes assumption

Challenge: $p(x_i | y, x_{<i})$ is hard to model (and estimate), especially for large i .

Solution:

Naive Bayes assumption

Features are **conditionally independent** given the label:

$$p(x | y) = \prod_{i=1}^d p(x_i | y). \quad (27)$$

A strong assumption in general, but works well in practice.

Parametrize $p(x_i | y)$ and $p(y)$

For binary x_i , assume $p(x_i | y)$ follows Bernoulli distributions.

$$p(x_i = 1 | y = 1) = \theta_{i,1}, \quad p(x_i = 1 | y = 0) = \theta_{i,0}. \quad (28)$$

Parametrize $p(x_i | y)$ and $p(y)$

For binary x_i , assume $p(x_i | y)$ follows Bernoulli distributions.

$$p(x_i = 1 | y = 1) = \theta_{i,1}, \quad p(x_i = 1 | y = 0) = \theta_{i,0}. \quad (28)$$

Similarly,

$$p(y = 1) = \theta_0. \quad (29)$$

Parametrize $p(x_i | y)$ and $p(y)$

For binary x_i , assume $p(x_i | y)$ follows Bernoulli distributions.

$$p(x_i = 1 | y = 1) = \theta_{i,1}, \quad p(x_i = 1 | y = 0) = \theta_{i,0}. \quad (28)$$

Similarly,

$$p(y = 1) = \theta_0. \quad (29)$$

Thus,

$$p(x, y) = p(x | y)p(y) \quad (30)$$

$$= p(y) \prod_{i=1}^d p(x_i | y) \quad \text{NB assumption} \quad (31)$$

$$= p(y) \prod_{i=1}^d \theta_{i,y} \mathbb{I}\{x_i = 1\} + (1 - \theta_{i,y}) \mathbb{I}\{x_i = 0\} \quad (32)$$

Indicator function $\mathbb{I}\{\text{condition}\}$ evaluates to 1 if “condition” is true and 0 otherwise.

MLE for our NB model

We maximize the likelihood of the data $\prod_{n=1}^N p_{\theta}(x^{(n)}, y^{(n)})$ (as opposed to the *conditional* likelihood we've seen before).

Set $\frac{\partial}{\partial \theta_{j,1}} \ell$ to zero:

$$\theta_{j,1} = \frac{\sum_{n=1}^N \mathbb{I}\{y^{(n)} = 1 \wedge x_j^{(n)} = 1\}}{\sum_{n=1}^N \mathbb{I}\{y^{(n)} = 1\}} \quad (33)$$

MLE for our NB model

We maximize the likelihood of the data $\prod_{n=1}^N p_{\theta}(x^{(n)}, y^{(n)})$ (as opposed to the *conditional* likelihood we've seen before).

Set $\frac{\partial}{\partial \theta_{j,1}} \ell$ to zero:

$$\theta_{j,1} = \frac{\sum_{n=1}^N \mathbb{I}\{y^{(n)} = 1 \wedge x_j^{(n)} = 1\}}{\sum_{n=1}^N \mathbb{I}\{y^{(n)} = 1\}} \quad (33)$$

In practice, count words:

$$\frac{\text{number of fake reviews containing "absolutely"}}{\text{number of fake reviews}}$$

Exercise: show that

$$\theta_{j,0} = \frac{\sum_{n=1}^N \mathbb{I}\{y^{(n)} = 0 \wedge x_j^{(n)} = 1\}}{\sum_{n=1}^N \mathbb{I}\{y^{(n)} = 0\}}, \theta_0 = \frac{\sum_{n=1}^N \mathbb{I}\{y^{(n)} = 1\}}{N} \quad (34)$$

NB assumption: **conditionally independent** features given the label

Recipe for learning a NB model:

- 1 Choose $p(x_i | y)$, e.g., Bernoulli distribution for binary x_i .
- 2 Choose $p(y)$, often a categorical distribution.
- 3 Estimate parameters by MLE (same as the strategy for conditional models) .

Next, NB with continuous features.

NB with continuous inputs

Let's consider a multiclass classification task with continuous inputs.

$$p(x_i | y) \sim \mathcal{N}(\mu_{i,y}, \sigma_{i,y}^2) \quad (35)$$

$$p(y = k) = \theta_k \quad (36)$$

NB with continuous inputs

Let's consider a multiclass classification task with continuous inputs.

$$p(x_i | y) \sim \mathcal{N}(\mu_{i,y}, \sigma_{i,y}^2) \quad (35)$$

$$p(y = k) = \theta_k \quad (36)$$

Likelihood of the data:

$$p(\mathcal{D}) = \prod_{n=1}^N p(y^{(n)}) \prod_{i=1}^d p(x_i^{(n)} | y^{(n)}) \quad (37)$$

$$= \prod_{n=1}^N \theta_{y^{(n)}} \prod_{i=1}^d \frac{1}{\sqrt{2\pi}\sigma_{i,y^{(n)}}} \exp\left(-\frac{1}{2\sigma_{i,y^{(n)}}^2} \left(x_i^{(n)} - \mu_{i,y^{(n)}}\right)^2\right) \quad (38)$$

MLE for Gaussian NB

Log likelihood:

$$\ell = \sum_{n=1}^N \log \theta_{y^{(n)}} + \sum_{n=1}^N \sum_{i=1}^d \log \frac{1}{\sqrt{2\pi}\sigma_{i,y^{(n)}}} - \frac{1}{2\sigma_{i,y^{(n)}}^2} \left(x_i^{(n)} - \mu_{i,y^{(n)}} \right)^2 \quad (39)$$

(40)

(41)

MLE for Gaussian NB

Log likelihood:

$$\ell = \sum_{n=1}^N \log \theta_{y^{(n)}} + \sum_{n=1}^N \sum_{i=1}^d \log \frac{1}{\sqrt{2\pi}\sigma_{i,y^{(n)}}} - \frac{1}{2\sigma_{i,y^{(n)}}^2} \left(x_i^{(n)} - \mu_{i,y^{(n)}} \right)^2 \quad (39)$$

(40)

Set $\frac{\partial}{\partial \mu_{j,k}} \ell$ to zero:

(41)

MLE for Gaussian NB

Log likelihood:

$$\ell = \sum_{n=1}^N \log \theta_{y^{(n)}} + \sum_{n=1}^N \sum_{i=1}^d \log \frac{1}{\sqrt{2\pi}\sigma_{i,y^{(n)}}} - \frac{1}{2\sigma_{i,y^{(n)}}^2} \left(x_i^{(n)} - \mu_{i,y^{(n)}} \right)^2 \quad (39)$$

(40)

Set $\frac{\partial}{\partial \mu_{j,k}} \ell$ to zero:

$$\mu_{j,k} = \frac{\sum_{n:y^{(n)}=k} x_j^{(n)}}{\sum_{n:y^{(n)}=k} 1} \quad (41)$$

MLE for Gaussian NB

Log likelihood:

$$\ell = \sum_{n=1}^N \log \theta_{y^{(n)}} + \sum_{n=1}^N \sum_{i=1}^d \log \frac{1}{\sqrt{2\pi}\sigma_{i,y^{(n)}}} - \frac{1}{2\sigma_{i,y^{(n)}}^2} \left(x_i^{(n)} - \mu_{i,y^{(n)}} \right)^2 \quad (39)$$

(40)

Set $\frac{\partial}{\partial \mu_{j,k}} \ell$ to zero:

$$\mu_{j,k} = \frac{\sum_{n:y^{(n)}=k} x_j^{(n)}}{\sum_{n:y^{(n)}=k} 1} = \text{sample mean of } x_j \text{ in class } k \quad (41)$$

Show that

$$\sigma_{j,k}^2 = \frac{\sum_{n:y^{(n)}=k} \left(x_j^{(n)} - \mu_{j,k}\right)^2}{\sum_{n:y^{(n)}=k} 1} = \text{sample variance of } x_j \text{ in class } k \quad (42)$$

$$\theta_k = \frac{\sum_{n:y^{(n)}=k} 1}{N} \quad (\text{class prior}) \quad (43)$$

Decision boundary of the Gaussian NB model

Is the Gaussian NB model a linear classifier?

Decision boundary of the Gaussian NB model

Is the Gaussian NB model a linear classifier?

$$\log \frac{p(y = 1 | x)}{p(y = 0 | x)} = \log \frac{p(x | y = 1)p(y = 1)}{p(x | y = 0)p(y = 0)} \quad (44)$$

(48)

Decision boundary of the Gaussian NB model

Is the Gaussian NB model a linear classifier?

$$\log \frac{p(y = 1 | x)}{p(y = 0 | x)} = \log \frac{p(x | y = 1)p(y = 1)}{p(x | y = 0)p(y = 0)} \quad (44)$$

$$= \log \frac{\theta_0}{1 - \theta_0} + \sum_{i=1}^d \left(\log \sqrt{\frac{\sigma_{i,0}^2}{\sigma_{i,1}^2}} + \left(\frac{(x_i - \mu_{i,0})^2}{2\sigma_{i,0}^2} - \frac{(x_i - \mu_{i,1})^2}{2\sigma_{i,1}^2} \right) \right)$$

(48)

Decision boundary of the Gaussian NB model

Is the Gaussian NB model a linear classifier?

$$\log \frac{p(y = 1 | x)}{p(y = 0 | x)} = \log \frac{p(x | y = 1)p(y = 1)}{p(x | y = 0)p(y = 0)} \quad (44)$$

$$= \log \frac{\theta_0}{1 - \theta_0} + \sum_{i=1}^d \left(\log \sqrt{\frac{\sigma_{i,0}^2}{\sigma_{i,1}^2}} + \left(\frac{(x_i - \mu_{i,0})^2}{2\sigma_{i,0}^2} - \frac{(x_i - \mu_{i,1})^2}{2\sigma_{i,1}^2} \right) \right) \quad \text{quadratic} \quad (45)$$

(48)

Decision boundary of the Gaussian NB model

Is the Gaussian NB model a linear classifier?

$$\log \frac{p(y=1|x)}{p(y=0|x)} = \log \frac{p(x|y=1)p(y=1)}{p(x|y=0)p(y=0)} \quad (44)$$

$$= \log \frac{\theta_0}{1-\theta_0} + \sum_{i=1}^d \left(\log \sqrt{\frac{\sigma_{i,0}^2}{\sigma_{i,1}^2}} + \left(\frac{(x_i - \mu_{i,0})^2}{2\sigma_{i,0}^2} - \frac{(x_i - \mu_{i,1})^2}{2\sigma_{i,1}^2} \right) \right) \quad \text{quadratic} \quad (45)$$

$$\text{assume that } \sigma_{i,0} = \sigma_{i,1} = \sigma_i, \quad (\theta_0 = 0.5) \quad (46)$$

$$= \sum_{i=1}^d \frac{1}{2\sigma_i^2} \left((x_i - \mu_{i,0})^2 - (x_i - \mu_{i,1})^2 \right) \quad (47)$$

$$= \sum_{i=1}^d \frac{\mu_{i,1} - \mu_{i,0}}{\sigma_i^2} x_i + \frac{\mu_{i,0}^2 - \mu_{i,1}^2}{2\sigma_i^2} \quad \text{linear} \quad (48)$$

Decision boundary of the Gaussian NB model

Assuming the variance of each feature is the same for both classes, we have

$$\log \frac{p(y=1|x)}{p(y=0|x)} = \sum_{i=1}^d \frac{\mu_{i,1} - \mu_{i,0}}{\sigma_i^2} x_i + \frac{\mu_{i,0}^2 - \mu_{i,1}^2}{2\sigma_i^2} \quad (49)$$

$$= \theta^T x \quad \text{where else have we seen it?} \quad (50)$$

$$(51)$$

Decision boundary of the Gaussian NB model

Assuming the variance of each feature is the same for both classes, we have

$$\log \frac{p(y=1|x)}{p(y=0|x)} = \sum_{i=1}^d \frac{\mu_{i,1} - \mu_{i,0}}{\sigma_i^2} x_i + \frac{\mu_{i,0}^2 - \mu_{i,1}^2}{2\sigma_i^2} \quad (49)$$

$$= \theta^T x$$

where else have we seen it? (50)

(51)

$$\theta_i = \frac{\mu_{i,1} - \mu_{i,0}}{\sigma_i^2} \quad \text{for } i \in [1, d] \quad (52)$$

$$\theta_0 = \sum_{i=1}^d \frac{\mu_{i,0}^2 - \mu_{i,1}^2}{2\sigma_i^2} \quad \text{bias term} \quad (53)$$

Naive Bayes vs logistic regression

	logistic regression	Gaussian naive Bayes
model type	conditional/discriminative	generative
parametrization	$p(y x)$	$p(x y), p(y)$
assumptions on Y	Bernoulli	Bernoulli
assumptions on X	—	Gaussian
decision boundary	$\theta_{\text{LR}}^T x$	$\theta_{\text{GNB}}^T x$

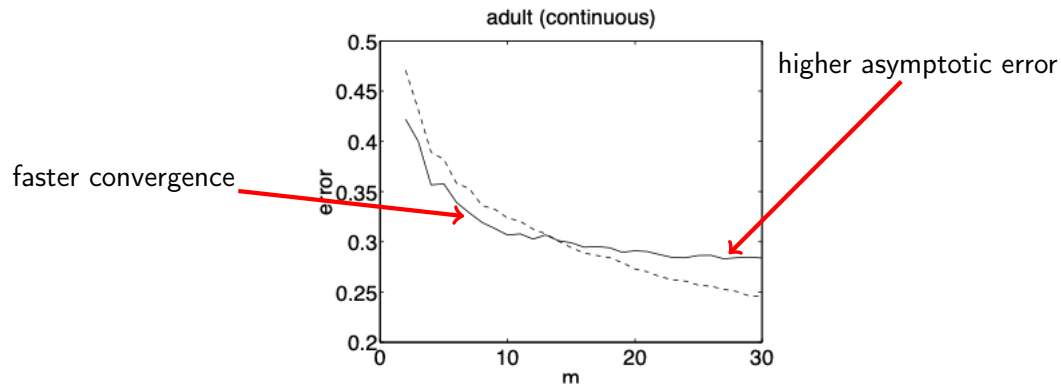
Naive Bayes vs logistic regression

	logistic regression	Gaussian naive Bayes
model type	conditional/discriminative	generative
parametrization	$p(y x)$	$p(x y), p(y)$
assumptions on Y	Bernoulli	Bernoulli
assumptions on X	—	Gaussian
decision boundary	$\theta_{\text{LR}}^T x$	$\theta_{\text{GNB}}^T x$

Given the same training data, is $\theta_{\text{LR}} = \theta_{\text{GNB}}$?

Generative vs discriminative classifiers

Ng, A. and Jordan, M. (2002). [On discriminative versus generative classifiers: A comparison of logistic regression and naive Bayes](#). In Advances in Neural Information Processing Systems 14.



Solid line: naive Bayes; dashed line: logistic regression.

Naive Bayes vs logistic regression

Logistic regression and Gaussian naive Bayes converge to the same classifier asymptotically, assuming the GNB assumption holds.

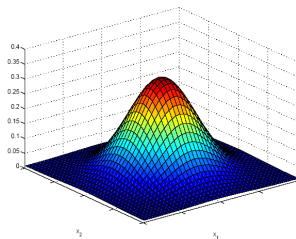
- Data points are generated from Gaussian distributions for each class
- Each dimension is independently generated
- Shared variance

What if the GNB assumption is not true?

Multivariate Gaussian Distribution

- $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, a Gaussian (or normal) distribution defined as

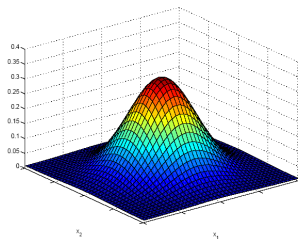
$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$



Multivariate Gaussian Distribution

- $x \sim \mathcal{N}(\mu, \Sigma)$, a Gaussian (or normal) distribution defined as

$$p(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right]$$

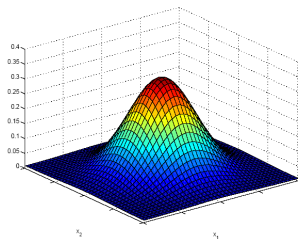


- Mahalanobis distance $(x - \mu_k)^T \Sigma^{-1} (x - \mu_k)$ measures the distance from x to μ in terms of Σ

Multivariate Gaussian Distribution

- $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, a Gaussian (or normal) distribution defined as

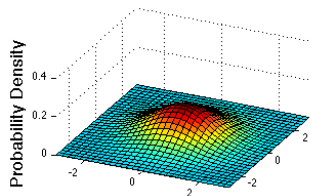
$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$



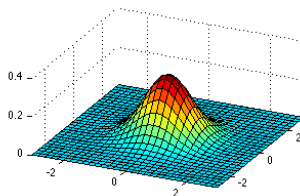
- Mahalanobis distance $(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)$ measures the distance from \mathbf{x} to $\boldsymbol{\mu}$ in terms of $\boldsymbol{\Sigma}$
- It normalizes for difference in variances and correlations

Bivariate Normal

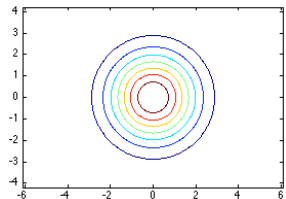
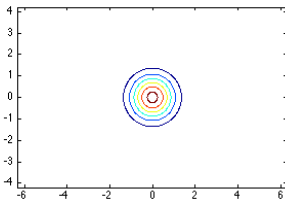
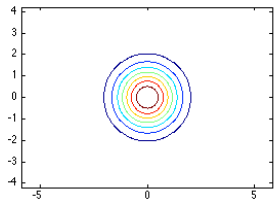
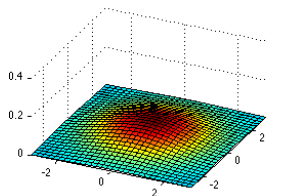
$$\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$



$$\Sigma = 0.5 \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

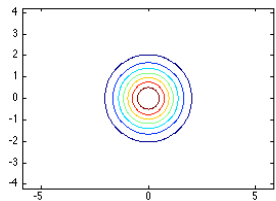
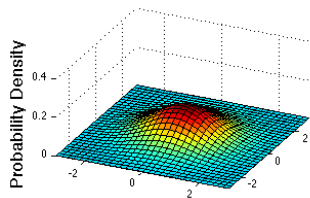


$$\Sigma = 2 \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

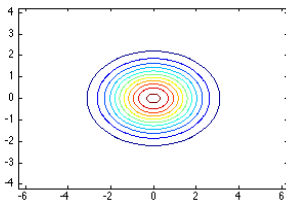
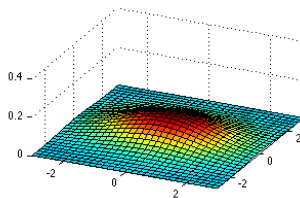


Bivariate Normal

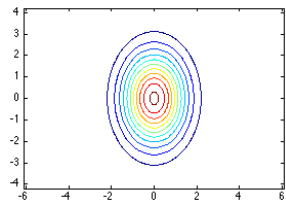
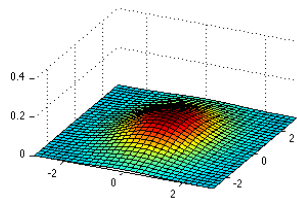
$$\text{var}(x_1) = \text{var}(x_2)$$



$$\text{var}(x_1) > \text{var}(x_2)$$

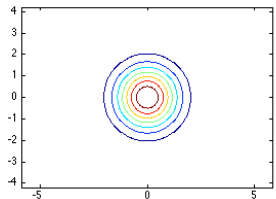
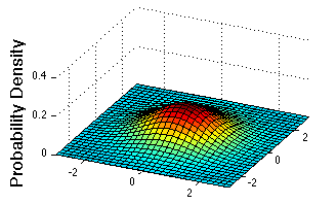


$$\text{var}(x_1) < \text{var}(x_2)$$

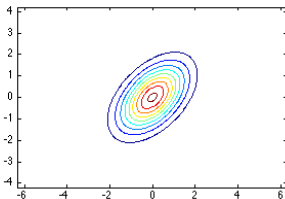
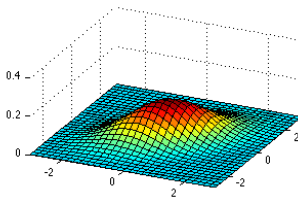


Bivariate Normal

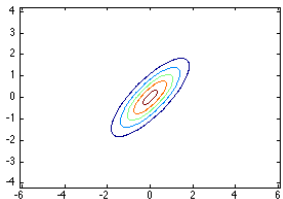
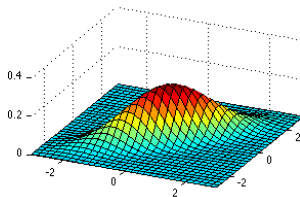
$$\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$



$$\Sigma = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}$$

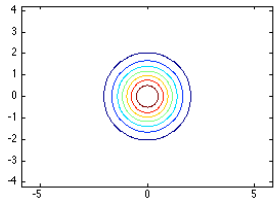
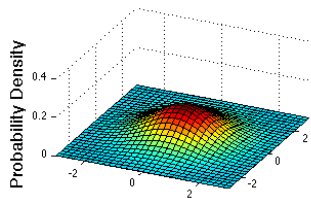


$$\Sigma = \begin{pmatrix} 1 & 0.8 \\ 0.8 & 1 \end{pmatrix}$$

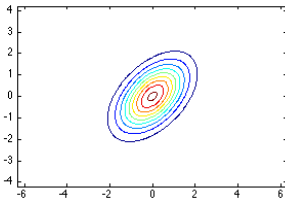
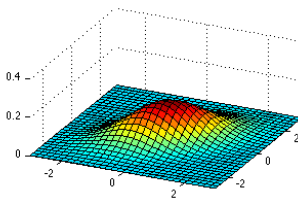


Bivariate Normal

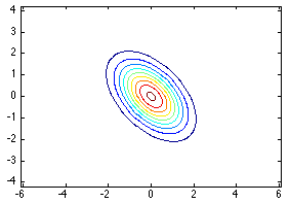
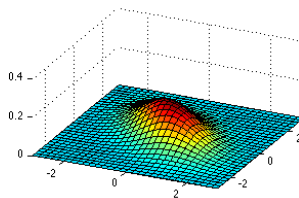
$$\text{Cov}(x_1, x_2) = 0$$



$$\text{Cov}(x_1, x_2) > 0$$



$$\text{Cov}(x_1, x_2) < 0$$



Gaussian Bayes Classifier

- Gaussian Bayes Classifier in its general form assumes that $p(\mathbf{x}|y)$ is distributed according to a multivariate normal (Gaussian) distribution
- Multivariate Gaussian distribution:

$$p(\mathbf{x}|t = k) = \frac{1}{(2\pi)^{d/2}|\Sigma_k|^{1/2}} \exp \left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k) \right]$$

where $|\Sigma_k|$ denotes the determinant of the matrix, and d is dimension of \mathbf{x}

Gaussian Bayes Classifier

- Gaussian Bayes Classifier in its general form assumes that $p(\mathbf{x}|\mathbf{y})$ is distributed according to a multivariate normal (Gaussian) distribution
- Multivariate Gaussian distribution:

$$p(\mathbf{x}|\mathbf{t} = k) = \frac{1}{(2\pi)^{d/2}|\Sigma_k|^{1/2}} \exp \left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k) \right]$$

where $|\Sigma_k|$ denotes the determinant of the matrix, and d is dimension of \mathbf{x}

- Each class k has associated mean vector $\boldsymbol{\mu}_k$ and covariance matrix Σ_k

Gaussian Bayes Classifier

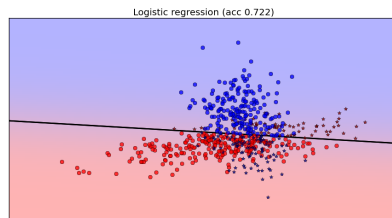
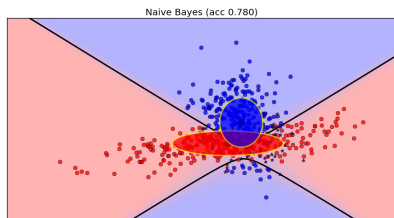
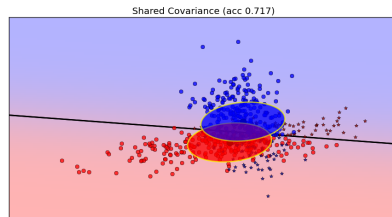
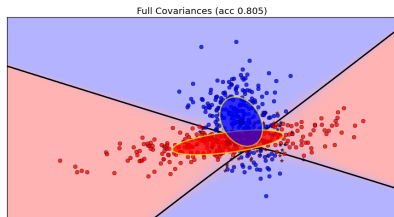
- Gaussian Bayes Classifier in its general form assumes that $p(\mathbf{x}|y)$ is distributed according to a multivariate normal (Gaussian) distribution
- Multivariate Gaussian distribution:

$$p(\mathbf{x}|t = k) = \frac{1}{(2\pi)^{d/2}|\Sigma_k|^{1/2}} \exp \left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k) \right]$$

where $|\Sigma_k|$ denotes the determinant of the matrix, and d is dimension of \mathbf{x}

- Each class k has associated mean vector $\boldsymbol{\mu}_k$ and covariance matrix Σ_k
- Σ_k has $\mathcal{O}(d^2)$ parameters - could be hard to estimate

Example



Gaussian Bayes Binary Classifier Cases

Different cases on the covariance matrix:

- Full covariance: Quadratic decision boundary
- Shared covariance: Linear decision boundary
- Naive Bayes: Diagonal covariance matrix, quadratic decision boundary

Gaussian Bayes Binary Classifier Cases

Different cases on the covariance matrix:

- Full covariance: Quadratic decision boundary
- Shared covariance: Linear decision boundary
- Naive Bayes: Diagonal covariance matrix, quadratic decision boundary

GBC vs. Logistic Regression:

- If data is truly Gaussian distributed, then shared covariance = logistic regression.
- But logistic regression can learn other distributions.

Summary

- Probabilistic framework of using maximum likelihood as a more principled way to derive loss functions.
- Conditional vs. generative
- Generative models the joint distribution, and may lead to more assumption on the data.

Summary

- Probabilistic framework of using maximum likelihood as a more principled way to derive loss functions.
- Conditional vs. generative
- Generative models the joint distribution, and may lead to more assumption on the data.
- When there is very few data point, it may be hard to model the distribution.
- Is there an equivalent “regularization” in a probabilistic framework?