

Classification



Yuling Yan

University of Wisconsin-Madison, Fall 2025

Classification problem

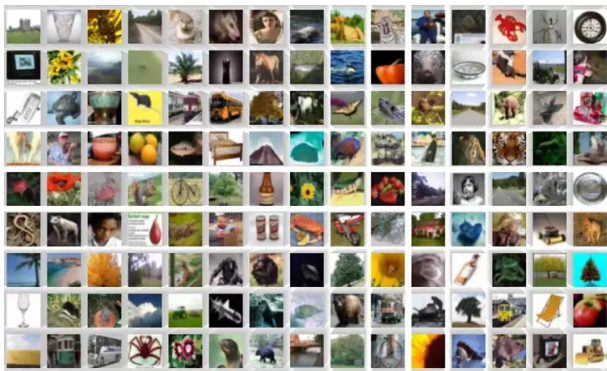
- Classification: assign a label (or category, class) to an observation based on its features
- \mathcal{X} : input space (e.g. \mathbb{R}^d); \mathcal{Y} : output space (e.g. $\{1, 2, \dots, K\}$)
- $x \in \mathcal{X}$: feature vector, input, data point...
- $y \in \mathcal{Y}$: label, category, class...
- Classifier: a mapping $f : \mathcal{X} \rightarrow \mathcal{Y}$
- Goal: construct a classifier f that accurately predicts the label y given the features x

MNIST dataset



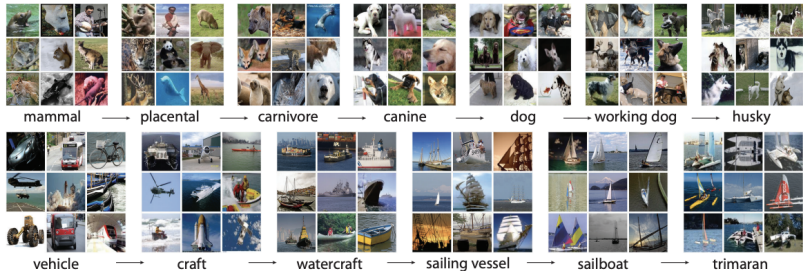
- Input: 28x28 gray scale (1 channel) images, i.e., $\mathcal{X} = \mathbb{R}^{28 \times 28}$ or \mathbb{R}^{784}
- Output: digits 0 through 9 (i.e., $\mathcal{Y} = \{0, 1, \dots, 9\}$)

CIFAR datasets



- Input: 32×32 RGB color (3 channels) images, i.e., $\mathcal{X} = \mathbb{R}^{32 \times 32 \times 3}$ or \mathbb{R}^{3072}
- Output: 10 classes (airplanes, cars, birds, cats, deer, dogs, frogs, horses, ships, and trucks) or 100 classes

ImageNet dataset



- Input: varies, often high-resolution (often $224 \times 224 \times 3$)
- Output: 1000 different categories

Mathematical set-up

- Modeling assumption: the data (input-output pairs) come from an underlying data distribution ρ over $\mathcal{X} \times \mathcal{Y}$
- Training data: $(x_1, y_1), \dots, (x_n, y_n) \stackrel{\text{i.i.d.}}{\sim} \rho$
- Error metric: for any given classifier f , its risk, defined as the average (expected) classification error on a new data is

$$R(f) := \mathbb{P}_{(X,Y) \sim \rho}(f(X) \neq Y)$$

- Supervised learning: build a classifier f based on training data, that makes the average classification error as small as possible

Questions

- Does there exist a “best” classifier?
— *this lecture*
- Can we construct this “best” classifier with the information of ρ ?
— *this and next lecture*
- What can we do when we only have a finite number of training data?
— *first half of the semester*

Bayes optimal classifier: binary case

- Consider the binary case: $\mathcal{Y} = \{0, 1\}$
- Define the Bayes classifier: for any $x \in \mathcal{X}$,

$$f^*(x) := \begin{cases} 1, & \text{if } \mathbb{P}(Y = 1 \mid X = x) \geq \mathbb{P}(Y = 0 \mid X = x), \\ 0, & \text{otherwise.} \end{cases}$$

Theorem 2.1 (Bayes optimal classifier: binary case)

The Bayes classifier f^ minimizes the misclassification error, i.e.,*

$$f^* \in \arg \min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{P}_{(X,Y) \sim \rho}(f(X) \neq Y).$$

Proof of Theorem 2.1

We need to show that, for any classifier $f : \mathcal{X} \rightarrow \mathcal{Y}$,

$$R(f) = \mathbb{P}(f(X) \neq Y) \geq \mathbb{P}(f^*(X) \neq Y) = R(f^*)$$

By tower property,

$$\begin{aligned}\mathbb{P}(f(X) \neq Y) &= \mathbb{E} [\mathbf{1}_{f(X) \neq Y}] \\&= \mathbb{E}_X [\mathbb{E} [\mathbf{1}_{f(X) \neq Y} \mid X]] && \text{(tower property)} \\&= \mathbb{E}_X [\mathbb{P}(f(X) \neq Y \mid X)] \\&\geq \mathbb{E}_X [\mathbb{P}(f^*(X) \neq Y \mid X)] && \text{(why?)} \\&= \mathbb{E}_X [\mathbb{E} [\mathbf{1}_{f^*(X) \neq Y} \mid X]] \\&= \mathbb{E} [\mathbf{1}_{f^*(X) \neq Y}] && \text{(tower property)} \\&= \mathbb{P}(f^*(X) \neq Y).\end{aligned}$$

It suffices to check

$$\mathbb{P}(f(X) \neq Y \mid X) \geq \mathbb{P}(f^*(X) \neq Y \mid X).$$

Proof of Theorem 2.1 (cont.)

Observe that

$$\begin{aligned}\mathbb{P}(f^*(X) \neq Y \mid X) &= \begin{cases} \mathbb{P}(Y = 0 \mid X) & \text{if } \mathbb{P}(Y = 1 \mid X) \geq \mathbb{P}(Y = 0 \mid X) \\ \mathbb{P}(Y = 1 \mid X) & \text{if } \mathbb{P}(Y = 1 \mid X) < \mathbb{P}(Y = 0 \mid X) \end{cases} \\ &= \min \{ \mathbb{P}(Y = 1 \mid X), \mathbb{P}(Y = 0 \mid X) \}\end{aligned}$$

and

$$\begin{aligned}\mathbb{P}(f(X) \neq Y \mid X) &= \begin{cases} \mathbb{P}(Y = 0 \mid X) & \text{if } f(X) = 1 \\ \mathbb{P}(Y = 1 \mid X) & \text{if } f(X) = 0 \end{cases} \\ &\geq \min \{ \mathbb{P}(Y = 1 \mid X), \mathbb{P}(Y = 0 \mid X) \}.\end{aligned}$$

Therefore

$$\mathbb{P}(f^*(X) \neq Y \mid X) \geq \mathbb{P}(f(X) \neq Y \mid X).$$

A few remarks

Bayes optimal classifier

$$f^*(x) := \begin{cases} 1, & \text{if } \mathbb{P}(Y = 1 \mid X = x) \geq \mathbb{P}(Y = 0 \mid X = x), \\ 0, & \text{otherwise.} \end{cases}$$

- Depends on the true underlying data distribution ρ
- The optimal classifier might not be unique
- When \mathcal{X} is discrete, it is equivalent to

$$f^*(x) := \begin{cases} 1, & \text{if } \mathbb{P}(X = x, Y = 1) \geq \mathbb{P}(X = x, Y = 0), \\ 0, & \text{otherwise.} \end{cases}$$

Bayes risk: binary case

- Bayes risk:

$$R^* := \mathbb{P}_{(X,Y) \sim \rho}(f^*(X) \neq Y)$$

- The Bayes risk serves as a lower bound for the classification error that any practical classifier can achieve:

$$R^* = \min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{P}_{(X,Y) \sim \rho}(f(X) \neq Y).$$

- It represents the inherent uncertainty in the classification problem due to overlapping distributions of the classes.
- Excess risk: $R(f) - R^*$

Bayes optimal classifier: multiclass setting

- Consider the multiclass case: $\mathcal{Y} = \{1, \dots, K\}$
- Define the Bayes classifier: for any $x \in \mathcal{X}$,

$$f^*(x) := \arg \max_{y \in \mathcal{Y}} \mathbb{P}(Y = y \mid X = x)$$

Theorem 2.2 (Bayes optimal classifier: multiclass case)

The Bayes classifier f^ minimizes the misclassification error, i.e.,*

$$f^* \in \arg \min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{P}_{(X,Y) \sim \rho}(f(X) \neq Y).$$

Bayes optimal classifier: multiclass setting

- Consider the multiclass case: $\mathcal{Y} = \{1, \dots, K\}$
- Define the Bayes classifier: for any $x \in \mathcal{X}$,

$$f^*(x) := \arg \max_{y \in \mathcal{Y}} \mathbb{P}(Y = y \mid X = x)$$

Theorem 2.2 (Bayes optimal classifier: multiclass case)

The Bayes classifier f^ minimizes the misclassification error, i.e.,*

$$f^* \in \arg \min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{P}_{(X,Y) \sim \rho}(f(X) \neq Y).$$

Proof: similar to Theorem 2.1, it suffices to check for any classifier f

$$\mathbb{P}(f(X) \neq Y \mid X) \geq \mathbb{P}(f^*(X) \neq Y \mid X).$$

More general loss function?

- Consider more general loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$
- Define the risk for a classifier $f : \mathcal{X} \rightarrow \mathcal{Y}$ as

$$R_\ell(f) := \mathbb{E}_{(X,Y) \sim \rho}[\ell(f(X), Y)]$$

- Example: with 0-1 loss $\ell(y, y') = \mathbb{1}\{y \neq y'\}$, we recover the average classification error

$$R(f) = \mathbb{P}_{(X,Y) \sim \rho}(f(X) \neq Y)$$

- Goal: find f that minimizes the risk $R_\ell(f)$ (the Bayes classifier might not be optimal...)

Question: Can you think of settings where other types of loss functions are more appropriate than the 0-1 loss?

Example: traffic signs



- $\mathcal{Y} = \{\text{stop sign}, 50 \text{ mph}, 40 \text{ mph}\}$.
- Predicting 50 mph when it is actually a stop sign is worse than predicting 40 mph when it is actually 50mph.
- 0-1 loss is not suitable here...

Example: traffic signs



- $\mathcal{Y} = \{\text{stop sign}, 50 \text{ mph}, 40 \text{ mph}\}$.
- Predicting 50 mph when it is actually a stop sign is worse than predicting 40 mph when it is actually 50mph.
- 0-1 loss is not suitable here...

We will discuss classification with general loss later if time permits

Supervised learning

- Go back to 0-1 loss
- In practice, we don't know ρ . It is in general impossible to compute the Bayes classifier f^\star
- Goal: build a classifier $f : \mathcal{X} \rightarrow \mathcal{Y}$ based on training data $(x_1, y_1), \dots, (x_n, y_n) \stackrel{\text{i.i.d.}}{\sim} \rho$
- Hope: achieve small excess risk $R(f) - R^\star$
- High-level framework:
 - Make some modeling assumptions on ρ
 - Design a good classifier f under this setup
 - For example, a good classifier may satisfy

$$R(f) - R^\star \leq h(n)$$

where $h(n)$ is a function of the sample size n describing the rate of convergence, e.g., $h(n) = O(1/n)$.

Linear Methods for Classification

Linear classifiers

- Linear classifiers: decision boundaries are linear hyperplanes

- Hyperplane $\mathcal{H}_{\beta, \beta_0} = \{\mathbf{x} \in \mathbb{R}^d : \langle \beta, \mathbf{x} \rangle + \beta_0 = 0\}$
- Half planes cut by $\mathcal{H}_{\beta, \beta_0}$:

$$\mathcal{H}_{\beta, \beta_0}^+ = \{\mathbf{x} \in \mathbb{R}^d : \langle \beta, \mathbf{x} \rangle + \beta_0 \geq 0\},$$

$$\mathcal{H}_{\beta, \beta_0}^- = \{\mathbf{x} \in \mathbb{R}^d : \langle \beta, \mathbf{x} \rangle + \beta_0 < 0\}.$$

- Example: in the binary case, the linear classifier has the form

$$f(\mathbf{x}) = \mathbb{1}\{\mathbf{x} \in \mathcal{H}_{\beta, \beta_0}^+\}$$

- Three approaches to learn a linear classifier from the data:
 - Linear discriminant analysis (LDA)
 - Logistic regression
 - Perceptrons and Support vector machines (SVMs)

Linear discriminant analysis (LDA)

- Model set-up: $\mathcal{X} = \mathbb{R}^d$, $\mathcal{Y} = \{1, \dots, K\}$. For $k = 1, \dots, K$,

$$\mathbb{P}(Y = k) = \omega_k, \quad X \mid Y = k \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma})$$

where $\omega_k \geq 0$, $\sum_{k=1}^K \omega_k = 1$, $\boldsymbol{\mu}_k \in \mathbb{R}^d$, $\boldsymbol{\Sigma} \in \mathbb{S}_+^d$

- The Bayes classifier under this setup: for any \boldsymbol{x} , compute

$$\delta_k(\boldsymbol{x}) := \underbrace{\boldsymbol{x}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k - \frac{1}{2} \boldsymbol{\mu}_k^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k + \log \omega_k}_{\propto \log \mathbb{P}(Y=k \mid X=\boldsymbol{x}) + \text{constant}}.$$

Let $f^*(\boldsymbol{x}) = \arg \max_{1 \leq k \leq K} \delta_k(\boldsymbol{x})$.

- Issue: model parameters are unknown...

Plug-in approach

- Plug-in approach: replace the unknown parameters with reliable estimates
- Suppose we have i.i.d. data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \stackrel{\text{i.i.d.}}{\sim} \rho$
- For each $1 \leq k \leq K$, let $n_k = \sum_{i=1}^n \mathbb{1}\{y_i = k\}$ and

$$\hat{\boldsymbol{\mu}}_k = \frac{1}{n_k} \sum_{i: y_i = k} \mathbf{x}_i, \quad \hat{\omega}_k = \frac{n_k}{n}$$

- Estimate the covariance matrix

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{N - \textcolor{red}{K}} \sum_{k=1}^K \sum_{i: y_i = k} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)^\top$$

- Replace $\boldsymbol{\mu}_k, \omega_k, \boldsymbol{\Sigma}$ with $\hat{\boldsymbol{\mu}}_k, \hat{\omega}_k, \hat{\boldsymbol{\Sigma}}$

$$\hat{\delta}_k(\mathbf{x}) := \underbrace{\mathbf{x}^\top \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_k - \frac{1}{2} \hat{\boldsymbol{\mu}}_k^\top \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_k}_{\text{linear in } \mathbf{x}} + \log \hat{\omega}_k.$$

Generalization

- Consider a more general set-up: for $k = 1, \dots, K$, assume

$$\mathbb{P}(Y = k) = \omega_k, \quad X \mid Y = k \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

where $\omega_k \geq 0$, $\sum_{k=1}^K \omega_k = 1$, $\mu_k \in \mathbb{R}^d$, $\Sigma_k \in \mathbb{S}_+^d$

- This setup will lead to the so-called quadratic discriminant analysis (QDA)
- Homework: derive QDA
 - What is the Bayes classifier under this setup?
 - How to derive a practical (data-driven) classifier?
 - Is this still a linear classifier?

Logistic regression

- Model set-up: $\mathcal{X} = \mathbb{R}^d$, $\mathcal{Y} = \{0, 1, \dots, K\}$. Let

$$\mathbb{P}(Y = k \mid \mathbf{x}) = \frac{\exp(\boldsymbol{\beta}_k^\top \mathbf{x} + \beta_{0,k})}{1 + \sum_{k'=1}^K \exp(\boldsymbol{\beta}_{k'}^\top \mathbf{x} + \beta_{0,k'})}, \quad (1 \leq k \leq K),$$
$$\mathbb{P}(Y = 0 \mid \mathbf{x}) = \frac{1}{1 + \sum_{k'=1}^K \exp(\boldsymbol{\beta}_{k'}^\top \mathbf{x} + \beta_{0,k'})},$$

where the parameters $\boldsymbol{\beta}_k \in \mathbb{R}^d$, $\beta_{0,k} \in \mathbb{R}$ for $k = 1, \dots, K$

Logistic regression

- Model set-up: $\mathcal{X} = \mathbb{R}^d \times \{\mathbf{1}\}$, $\mathcal{Y} = \{0, 1, \dots, K\}$. Let

$$\mathbb{P}(Y = k \mid \mathbf{x}) = \frac{\exp(\boldsymbol{\beta}_k^\top \mathbf{x})}{1 + \sum_{k'=1}^K \exp(\boldsymbol{\beta}_{k'}^\top \mathbf{x})}, \quad (k = 1, \dots, K),$$
$$\mathbb{P}(Y = 0 \mid \mathbf{x}) = \frac{1}{1 + \sum_{k'=1}^K \exp(\boldsymbol{\beta}_{k'}^\top \mathbf{x})},$$

where the parameters $\boldsymbol{\beta}_k \in \mathbb{R}^{d+1}$ for $k = 1, \dots, K$

- Bayes classifier:

$$f(\mathbf{x}) = \begin{cases} \operatorname{argmax}_{1 \leq k \leq K} \boldsymbol{\beta}_k^\top \mathbf{x}, & \text{if } \max_{1 \leq k \leq K} \boldsymbol{\beta}_k^\top \mathbf{x} > 0, \\ 0, & \text{otherwise.} \end{cases}$$

- Estimate $\boldsymbol{\beta}_k$'s: maximum likelihood estimation (MLE)

Maximum likelihood estimation

- Suppose we have i.i.d. data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$
- The negative log-likelihood function

$$\ell(\boldsymbol{\beta}) = -\frac{1}{n} \sum_{k=1}^K \sum_{i: y_i=k} \mathbf{x}_i^\top \boldsymbol{\beta}_k + \frac{1}{n} \sum_{i=1}^n \log \left[1 + \sum_{k'=1}^K \exp(\mathbf{x}_i^\top \boldsymbol{\beta}_{k'}) \right]$$

- Maximum likelihood estimation (MLE)

$$\hat{\boldsymbol{\beta}} := \arg \min_{\boldsymbol{\beta}} \ell(\boldsymbol{\beta})$$

- Convex optimization: solve by e.g., gradient descent

$$\boldsymbol{\beta}^{t+1} = \boldsymbol{\beta}^t - \eta \nabla \ell(\boldsymbol{\beta}^t) \quad (t = 0, 1, \dots)$$

A brief introduction to gradient descent

Gradient descent (GD) for solving $\min_{\beta \in \mathbb{R}^d} L(\beta)$:

$$\beta^{t+1} = \beta^t - \eta \nabla L(\beta^t) \quad (t = 0, 1, \dots)$$

When η is properly small, GD satisfy the following properties:

- For smooth function L , GD is a descent algorithm: $L(\beta^{t+1}) \leq L(\beta^t)$
- For convex + smooth function L , GD satisfies

$$L(\beta^t) - L(\beta^*) \leq O\left(\frac{\|\beta^0 - \beta^*\|_2^2}{t}\right) \quad (t = 0, 1, \dots)$$

for any minimizer β^*

- For strongly convex + smooth function L , GD satisfies

$$\|\beta^{t+1} - \beta^*\|_2 \leq (1 - \kappa)^t \|\beta^0 - \beta^*\|_2 \quad (t = 0, 1, \dots)$$

for some $\kappa \in (0, 1)$, where β^* is the unique minimizer

Stochastic gradient descent

Consider the following empirical risk minimization problem

$$\min_{\beta \in \mathbb{R}^d} L(\beta) := \frac{1}{n} \sum_{i=1}^n g(\beta; \mathbf{x}_i),$$

where $\mathbf{x}_1, \dots, \mathbf{x}_n$ are training data points.

- **Stochastic gradient descent:** for $t = 0, 1, \dots$,

$$\beta^{t+1} = \beta^t - \eta \nabla g(\beta^t; \mathbf{x}_{i_t}) \quad \text{where} \quad \mathbf{x}_{i_t} \stackrel{\text{ind.}}{\sim} \text{Unif}\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

- **Gradient descent:** for $t = 0, 1, \dots$,

$$\beta^{t+1} = \beta^t - \eta \nabla L(\beta^t) = \beta^t - \eta \frac{1}{n} \sum_{i=1}^n \nabla g(\beta; \mathbf{x}_i)$$

- **Advantage of SGD:** much faster updates, especially for large datasets, but still enjoys nice properties (sometimes even better than GD!)

Gradient descent methods

Example: GD / SGD for logistic regression

Take-away: (stochastic) gradient descent is the default method for solving unconstrained optimization problem

— simple and effective!

Recommended reading materials: Lecture 1 and 10 of the course

Large-Scale Optimization for Data Science

by Prof. Yuxin Chen (UPenn); Lecture on GD and SGD

Perceptrons and SVMs

Linearly separable data

- Consider binary classification: $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{Y} = \{1, -1\}$
- Training data: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$
- **Linearly separable data:** \exists a separating hyperplane $\mathcal{H}_{\boldsymbol{\beta}, \beta_0}$ s.t.

$$y_i \cdot (\mathbf{x}_i^\top \boldsymbol{\beta} + \beta_0) > 0 \quad (i = 1, \dots, n)$$

Linearly separable data

- Consider binary classification: $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{Y} = \{1, -1\}$
- Training data: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$
- **Linearly separable data:** \exists a separating hyperplane $\mathcal{H}_{\boldsymbol{\beta}, \beta_0}$ s.t.

$$y_i \cdot (\mathbf{x}_i^\top \boldsymbol{\beta} + \beta_0) > 0 \quad (i = 1, \dots, n)$$

- by merging β_0 into $\boldsymbol{\beta}$ and adding 1 to \mathbf{x}_i 's, this assumption becomes: $\exists \boldsymbol{\beta}_{\text{sep}} \in \mathbb{R}^{d+1}$

$$y_i \cdot \mathbf{x}_i^\top \boldsymbol{\beta}_{\text{sep}} > 0 \quad (i = 1, \dots, n)$$

Linearly separable data

- Consider binary classification: $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{Y} = \{1, -1\}$
- Training data: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$
- **Linearly separable data:** \exists a separating hyperplane $\mathcal{H}_{\boldsymbol{\beta}, \beta_0}$ s.t.

$$y_i \cdot (\mathbf{x}_i^\top \boldsymbol{\beta} + \beta_0) > 0 \quad (i = 1, \dots, n)$$

- by merging β_0 into $\boldsymbol{\beta}$ and adding 1 to \mathbf{x}_i 's, this assumption becomes: $\exists \boldsymbol{\beta}_{\text{sep}} \in \mathbb{R}^{d+1}$

$$y_i \cdot \mathbf{x}_i^\top \boldsymbol{\beta}_{\text{sep}} > 0 \quad (i = 1, \dots, n)$$

- **Goal:** search a separating hyperplane indexed by $\hat{\boldsymbol{\beta}}$

$$y_i \cdot \mathbf{x}_i^\top \hat{\boldsymbol{\beta}} > 0 \quad (i = 1, \dots, n)$$

(note that $\boldsymbol{\beta}_{\text{sep}}$ is not known a priori)

Perceptron Learning Algorithm

- For every $\beta \in \mathbb{R}^{d+1}$, define the set $\mathcal{M}_\beta := \underbrace{\{i : y_i \cdot \mathbf{x}_i^\top \beta \leq 0\}}_{\text{misclassified points}}$
- Target: minimize the perceptron loss

$$\sigma(\beta) := - \sum_{i \in \mathcal{M}_\beta} y_i \cdot \mathbf{x}_i^\top \beta \propto \sum_{i \in \mathcal{M}_\beta} \text{dist}(\mathbf{x}_i, \mathcal{H}_\beta)$$

where $\mathcal{H}_\beta = \{\mathbf{x} : \mathbf{x}^\top \beta = 0\}$

- Algorithm: initialize with $\beta^0 \in \mathbb{R}^{d+1}$, for $t = 0, 1, \dots$, update

$$\beta^{t+1} = \beta^t + \eta y_i \mathbf{x}_i, \quad \text{for a random } i \in \mathcal{M}_{\beta^t}$$

where $\eta > 0$ is the step size; in fact, we can take $\eta = 1$ here...

Perceptron Learning Algorithm

- For every $\beta \in \mathbb{R}^{d+1}$, define the set $\mathcal{M}_\beta := \underbrace{\{i : y_i \cdot \mathbf{x}_i^\top \beta \leq 0\}}_{\text{misclassified points}}$
- Target: minimize the perceptron loss

$$\sigma(\beta) := - \sum_{i \in \mathcal{M}_\beta} y_i \cdot \mathbf{x}_i^\top \beta \propto \sum_{i \in \mathcal{M}_\beta} \text{dist}(\mathbf{x}_i, \mathcal{H}_\beta)$$

where $\mathcal{H}_\beta = \{\mathbf{x} : \mathbf{x}^\top \beta = 0\}$

- Algorithm: initialize with $\beta^0 \in \mathbb{R}^{d+1}$, for $t = 0, 1, \dots$, update

$$\beta^{t+1} = \beta^t + y_i \mathbf{x}_i, \quad \text{for a random } i \in \mathcal{M}_{\beta^t}$$

- Interpretation: SGD with step size 1 (kind of...)

Convergence theory

Theorem 2.3

*When the data is **linearly separable**, the perceptron learning algorithm converges to a separating hyperplane in a finite number of steps.*

Convergence theory

Theorem 2.3

*When the data is **linearly separable**, the perceptron learning algorithm converges to a separating hyperplane in a finite number of steps.*

Limitations:

- solutions not unique: might converge to an **unstable** hyperplane

Convergence theory

Theorem 2.3

*When the data is **linearly separable**, the perceptron learning algorithm converges to a separating hyperplane in a finite number of steps.*

Limitations:

- solutions not unique: might converge to an **unstable** hyperplane
— resort to “**optimal separating hyperplane**”

Convergence theory

Theorem 2.3

When the data is linearly separable, the perceptron learning algorithm converges to a separating hyperplane in a finite number of steps.

Limitations:

- solutions not unique: might converge to an unstable hyperplane
— resort to “optimal separating hyperplane”
- only works linearly separable data. If the classes cannot be separated by a hyperplane, the algorithm will not converge

Convergence theory

Theorem 2.3

*When the data is **linearly separable**, the perceptron learning algorithm converges to a separating hyperplane in a finite number of steps.*

Limitations:

- solutions not unique: might converge to an **unstable** hyperplane
— resort to “**optimal separating hyperplane**”
- only works linearly separable data. If the classes cannot be separated by a hyperplane, the algorithm will not converge
- the “finite” number of steps can be very large

Optimal separating hyperplane

From now on, we “unmerge” β_0 from β , as they play different roles. Consider the optimization problem

$$\max_{\|\beta\|_2=1, \beta_0, M} M \quad \text{s.t.} \quad y_i(\mathbf{x}_i^\top \beta + \beta_0) \geq M \quad (i = 1, \dots, n)$$

Optimal separating hyperplane

From now on, we “unmerge” β_0 from β , as they play different roles.
Consider the optimization problem

$$\max_{\|\beta\|_2=1, \beta_0, M} M \quad \text{s.t.} \quad y_i(\mathbf{x}_i^\top \beta + \beta_0) \geq M \quad (i = 1, \dots, n)$$

Implications:

- the distance between \mathbf{x} and the hyperplane $\mathcal{H}_{\beta, \beta_0}$ is

$$\text{dist}(\mathbf{x}, \mathcal{H}_{\beta, \beta_0}) = \frac{|\beta^\top \mathbf{x} + \beta_0|}{\|\beta\|_2} \quad \text{if } \underline{\underline{\|\beta\|_2=1}} \quad |\beta^\top \mathbf{x} + \beta_0|$$

Optimal separating hyperplane

From now on, we “unmerge” β_0 from β , as they play different roles.
Consider the optimization problem

$$\max_{\|\beta\|_2=1, \beta_0, M} M \quad \text{s.t.} \quad y_i(\mathbf{x}_i^\top \beta + \beta_0) \geq M \quad (i = 1, \dots, n)$$

Implications:

- the distance between \mathbf{x} and the hyperplane $\mathcal{H}_{\beta, \beta_0}$ is $|\beta^\top \mathbf{x} + \beta_0|$
- offers a unique solution that maximizes the *margin* M
- **Margin:** the distance between $\mathcal{H}_{\beta, \beta_0}$ and the closest data points from each class
support vectors

Optimal separating hyperplane

From now on, we “unmerge” β_0 from β , as they play different roles. Consider the optimization problem

$$\max_{\|\beta\|_2=1, \beta_0, M} M \quad \text{s.t.} \quad y_i(\mathbf{x}_i^\top \beta + \beta_0) \geq M \quad (i = 1, \dots, n)$$

Implications:

- the distance between \mathbf{x} and the hyperplane $\mathcal{H}_{\beta, \beta_0}$ is $|\beta^\top \mathbf{x} + \beta_0|$
- offers a unique solution that maximizes the *margin* M
- **Margin:** the distance between $\mathcal{H}_{\beta, \beta_0}$ and the closest data points from each class
support vectors
- **Intuition:** a large margin on the training data will lead to good separation on the test data.

Reformulation as convex optimization

- **Original problem:**

$$\max_{\|\beta\|_2=1, \beta_0, M} M \quad \text{s.t.} \quad y_i(\mathbf{x}_i^\top \beta + \beta_0) \geq M \quad (i = 1, \dots, n)$$

- **Issue:** this is not a convex optimization problem...

Reformulation as convex optimization

- **Original problem:**

$$\max_{\|\beta\|_2=1, \beta_0, M} M \quad \text{s.t.} \quad y_i(\mathbf{x}_i^\top \beta + \beta_0) \geq M \quad (i = 1, \dots, n)$$

- **Issue:** this is not a convex optimization problem...
- **Reformulation:**

$$\min_{\beta, \beta_0} \|\beta\|_2^2 \quad \text{s.t.} \quad y_i(\mathbf{x}_i^\top \beta + \beta_0) \geq 1 \quad (i = 1, \dots, n)$$

this is a convex optimization problem

Reformulation as convex optimization

- **Original problem:**

$$\max_{\|\beta\|_2=1, \beta_0, M} M \quad \text{s.t.} \quad y_i(\mathbf{x}_i^\top \beta + \beta_0) \geq M \quad (i = 1, \dots, n)$$

- **Issue:** this is not a convex optimization problem...
- **Reformulation:**

$$\min_{\beta, \beta_0} \|\beta\|_2^2 \quad \text{s.t.} \quad y_i(\mathbf{x}_i^\top \beta + \beta_0) \geq 1 \quad (i = 1, \dots, n)$$

this is a convex optimization problem

- This is known as the support vector machine (SVM)

SVMs for separable data

$$\min_{\beta, \beta_0} \quad \frac{1}{2} \|\beta\|_2^2 \quad \text{s.t.} \quad y_i(\mathbf{x}_i^\top \beta + \beta_0) \geq 1 \quad (i = 1, \dots, n)$$

- SVM is a powerful method for binary classification
- finds a linear classifier with decision boundary $\{\mathbf{x} : \mathbf{x}^\top \hat{\beta} + \hat{\beta}_0 = 0\}$ to separate two classes with the maximum margin
- This is only feasible for *linearly separated data*
 - *can be generalized to accommodate non-separable data*
- What can we say about SVM?
 - *resort to duality theory!*

Convex optimization and duality theory

Primal problem and Lagrangian function

- Consider a convex optimization problem:

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) \quad \text{s.t.} \quad g_i(\mathbf{x}) \leq 0 \quad (i = 1, \dots, m).$$

where $f(\mathbf{x})$ and $g_i(\mathbf{x})$ are convex functions

- This is called the **primal problem**
- To handle the constraints, we introduce **Lagrange multipliers** λ_i
- The Lagrangian function is:

$$L(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x})$$

- What is the benefit of introducing the Lagrangian function?

The Dual Problem

Key observation:

$$\underbrace{\min_{\mathbf{x}: g(\mathbf{x}) \leq 0} f(\mathbf{x})}_{\text{primal problem}} \stackrel{(i)}{=} \min_{\mathbf{x}} \max_{\boldsymbol{\lambda} \geq 0} L(\mathbf{x}, \boldsymbol{\lambda}) \geq \max_{\boldsymbol{\lambda} \geq 0} \underbrace{\min_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda})}_{=: d(\boldsymbol{\lambda})} = \underbrace{\max_{\boldsymbol{\lambda} \geq 0} d(\boldsymbol{\lambda})}_{\text{dual problem}}$$

- relation (i) and (ii) always holds (why?)
- relation (ii) is often an equality (strong duality theory)
- The dual function $d(\boldsymbol{\lambda}) = \min_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda})$
- The **dual problem** is to maximize the dual function $d(\boldsymbol{\lambda})$:

$$\max_{\boldsymbol{\lambda} \geq 0} d(\boldsymbol{\lambda})$$

Strong and Weak Duality

Weak Duality: For any x feasible in the primal and any $\lambda \geq 0$, we have:

$$d(\lambda) \leq f(x)$$

Strong Duality: If the problem satisfies certain conditions (e.g., Slater's condition), then:

$$\min_{x: g(x) \leq 0} f(x) = \max_{\lambda \geq 0} d(\lambda)$$

- Slater's condition: the feasible region has an interior point, i.e.,

$$\exists x_0 \in \mathbb{R}^d \quad \text{s.t.} \quad g_i(x_0) < 0 \quad (i = 1, \dots, m).$$

- In convex optimization, **strong duality often holds**, meaning the primal and dual problems have the same optimal value.

KKT Conditions

The **Karush-Kuhn-Tucker (KKT)** conditions: if strong duality holds, and $(\mathbf{x}, \boldsymbol{\lambda})$ is the optimal solution pair for the **primal**/**dual** problem

$$\underbrace{\min_{\mathbf{x}: g(\mathbf{x}) \leq 0} f(\mathbf{x})}_{\text{primal problem}} = \underbrace{\max_{\boldsymbol{\lambda} \geq 0} d(\boldsymbol{\lambda})}_{\text{dual problem}},$$

then

- **Primal feasibility:** $g_i(\mathbf{x}) \leq 0$
- **Dual feasibility:** $\lambda_i \geq 0$
- **Complementary slackness:** $\lambda_i g_i(\mathbf{x}) = 0$
- **Stationarity:** $\nabla f(\mathbf{x}) + \sum_{i=1}^m \lambda_i \nabla g_i(\mathbf{x}) = 0$

— This is a necessary condition!

Back to SVMs

$$\min_{\beta, \beta_0} \quad \frac{1}{2} \|\beta\|_2^2 \quad \text{s.t.} \quad y_i(\mathbf{x}_i^\top \beta + \beta_0) \geq 1 \quad (i = 1, \dots, n)$$

- The **dual problem** for SVM is (why?):

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \quad \text{s.t.} \quad \sum_{i=1}^n \alpha_i y_i = 0, \alpha_i \geq 0$$

- It is straightforward to check that Slater's condition holds
— *primal and dual problems are equivalent!*
- The dual problem is a quadratic programming problem, which is easier to compute with standard software (e.g. CVX)

Checking KKT conditions

$$(P) \quad \min_{\beta, \beta_0} \quad \frac{1}{2} \|\beta\|_2^2 \quad \text{s.t.} \quad y_i(\mathbf{x}_i^\top \beta + \beta_0) \geq 1 \quad (i = 1, \dots, n)$$

$$(D) \quad \max_{\alpha} \quad \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \quad \text{s.t.} \quad \sum_{i=1}^n \alpha_i y_i = 0, \alpha_i \geq 0$$

The Karush-Kuhn-Tucker (KKT) conditions for optimality:

- **Primal feasibility:** $y_i(\beta^\top \mathbf{x}_i + \beta_0) \geq 1$
- **Dual feasibility:** $\alpha_i \geq 0$
- **Complementary slackness:** $\alpha_i[y_i(\beta^\top \mathbf{x}_i + \beta_0) - 1] = 0$
- **Stationarity:** $\beta = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$

Implications

For any optimal solution pair $(\beta^*, \beta_0^*, \alpha^*)$:

- **Support vectors:** data points \mathbf{x}_i with $\alpha_i > 0$

$$y_i(\beta^{*\top} \mathbf{x}_i + \beta_0^*) > 1 \implies \alpha_i = 0$$

$$\alpha_i > 0 \implies y_i(\beta^{*\top} \mathbf{x}_i + \beta_0^*) = 1$$

- **Recovering the primal solution:** after solving the dual problem (i.e., finding α_i^*), we can recover the primal solution (β^*, β_0^*) by

$$\beta^* = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i$$

and $\beta_0^* = y_i - \beta^{*\top} \mathbf{x}_i$ for any support vector \mathbf{x}_i

— β^* is a linear combination of the support vectors

Accommodating non-separable data

SVM for linearly separable data:

$$\min_{\beta, \beta_0} \quad \frac{1}{2} \|\beta\|_2^2 \quad \text{s.t.} \quad y_i(\mathbf{x}_i^\top \beta + \beta_0) \geq 1 \quad (i = 1, \dots, n)$$

- For non-separable data, we introduce slack variables $\xi_i \geq 0$ to allow violations of the margin:

$$\begin{aligned} \min_{\beta, \beta_0, \xi} \quad & \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i(\beta^\top \mathbf{x}_i + \beta_0) \geq 1 - \xi_i, \quad \xi_i \geq 0 \quad (i = 1, \dots, n) \end{aligned}$$

- $C > 0$ is the “cost” parameter
- the separable case corresponds to $C = \infty$

Dual problem: non-separable data

- **Primal problem:**

$$\begin{aligned} \min_{\beta, \beta_0, \xi} \quad & \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i(\beta^\top \mathbf{x}_i + \beta_0) \geq 1 - \xi_i, \quad \xi_i \geq 0 \quad (i = 1, \dots, n) \end{aligned}$$

- **Dual problem:**

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C \quad (i = 1, \dots, n) \end{aligned}$$

- Homework: derive the **dual problem** from the **primal problem**

Dual problem: non-separable data

- **Primal problem:**

$$\begin{aligned} \min_{\beta, \beta_0, \xi} \quad & \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i(\beta^\top \mathbf{x}_i + \beta_0) \geq 1 - \xi_i, \quad \xi_i \geq 0 \quad (i = 1, \dots, n) \end{aligned}$$

- **Dual problem:**

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C \quad (i = 1, \dots, n) \end{aligned}$$

- Homework: derive the **dual problem** from the **primal problem**

Kernel density classifier and naive Bayes classifier

Recap: Bayes optimal classifier

Bayes optimal classifier: for any $x \in \mathcal{X}$, output

$$f^*(x) := \arg \max_{y \in \mathcal{Y}} \mathbb{P}(Y = y \mid X = x)$$

- Issue: depends on unknown data distribution ρ

Recap: Bayes optimal classifier

Bayes optimal classifier: for any $x \in \mathcal{X}$, output

$$f^*(x) := \arg \max_{y \in \mathcal{Y}} \mathbb{P}(Y = y \mid X = x)$$

- Issue: depends on unknown data distribution ρ
- Bayes formula:

$$\mathbb{P}(Y = y \mid X = x) = \frac{\mathbb{P}(X = x \mid Y = y) \mathbb{P}(Y = y)}{\sum_{y' \in \mathcal{Y}} \mathbb{P}(X = x \mid Y = y') \mathbb{P}(Y = y')}$$

— *Is it possible to estimate these quantities?*

Recap: Bayes optimal classifier

Bayes optimal classifier: for any $x \in \mathcal{X}$, output

$$f^*(x) := \arg \max_{y \in \mathcal{Y}} \mathbb{P}(Y = y \mid X = x)$$

- Issue: depends on unknown data distribution ρ
- Bayes formula:

$$\mathbb{P}(Y = y \mid X = x) = \frac{\hat{\mathbb{P}}(X = x \mid Y = y) \hat{\mathbb{P}}(Y = y)}{\sum_{y' \in \mathcal{Y}} \hat{\mathbb{P}}(X = x \mid Y = y') \hat{\mathbb{P}}(Y = y')}$$

— *Is it possible to estimate these quantities?*

- Plug-in method:
 - marginal probabilities $\mathbb{P}(Y = y)$ are easy to estimate (use frequency)
 - key difficulty: **estimate conditional densities** $\mathbb{P}(X = x \mid Y = y)$

Detour: density estimation

Setup: density estimation

- Target: an unknown density function f
- What we have: i.i.d. data $X_1, \dots, X_n \sim f$
- Goal: construct a **good** density estimation $\hat{f}(\cdot)$ that satisfy

$$\hat{f}(x) \geq 0 \quad \text{and} \quad \int_0^1 \hat{f}(x) dx = 1$$

Setup: density estimation

- Target: an unknown density function f
- What we have: i.i.d. data $X_1, \dots, X_n \sim f$
- Goal: construct a **good** density estimation $\hat{f}(\cdot)$ that satisfy

$$\hat{f}(x) \geq 0 \quad \text{and} \quad \int_0^1 \hat{f}(x) dx = 1$$

- Criteria: mean integrated squared error (MISE)

$$\text{MISE}(\hat{f}) = \mathbb{E} \left[\int (\hat{f}(x) - f(x))^2 dx \right]$$

Setup: density estimation

- Target: an unknown density function f
- What we have: i.i.d. data $X_1, \dots, X_n \sim f$
- Goal: construct a **good** density estimation $\hat{f}(\cdot)$ that satisfy

$$\hat{f}(x) \geq 0 \quad \text{and} \quad \int_0^1 \hat{f}(x) dx = 1$$

- Criteria: mean integrated squared error (MISE)

$$\text{MISE}(\hat{f}) = \mathbb{E} \left[\int (\hat{f}(x) - f(x))^2 dx \right]$$

- Density estimation: find \hat{f} with as small MISE as possible
 - Histogram method
 - Kernel density estimation

Bias-variance tradeoff

Mean integrated squared error (MISE):

$$\text{MISE}(\hat{f}) = \mathbb{E} \left[\int (\hat{f}(x) - f(x))^2 dx \right]$$

- **Bias:** Measures how far the estimated density is from the true density on average.

$$b(x) := \mathbb{E}[\hat{f}(x)] - f(x)$$

- **Variance:** Measures how much $\hat{f}(x)$ fluctuates around its mean:

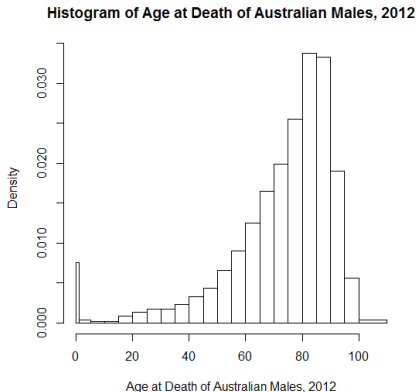
$$v(x) := \text{var}(\hat{f}(x)) = \mathbb{E}[(\hat{f}(x) - \mathbb{E}[\hat{f}(x)])^2]$$

Theorem 2.4

$$\text{MISE}(\hat{f}) = \int b^2(x) dx + \int v(x) dx$$

A simple approach: histogram

Histogram method: estimate the density by partitioning the interval and counting the frequency of data points in each partition



— credit to R.J. Oosterbaan

Histograms

- Consider 1D setting, and assume that $f(\cdot)$ is supported on $[0, 1]$
 - *we can always rescale the data to $[0, 1]$*
- The data is divided into m bins of equal width $h = 1/m$ (bandwidth)

$$B_1 = \left[0, \frac{1}{m}\right), \quad B_2 = \left[\frac{1}{m}, \frac{2}{m}\right), \quad \dots, \quad B_m = \left[\frac{m-1}{m}, 1\right]$$

- Each bin is assigned a probability proportional to the number of observations falling into that bin:

$$\hat{f}(x) := \begin{cases} \hat{p}_1/h, & x \in B_1, \\ \vdots & \vdots \\ \hat{p}_m/h, & x \in B_m, \end{cases} \quad \text{where} \quad \hat{p}_j = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \in B_j\}.$$

Optimal bandwidth

Theorem 2.5 (informal)

Under some regularity conditions, we have

$$\text{MISE}(\hat{f}) \approx \frac{h^2}{12} \int f'(u)^2 du + \frac{1}{nh}$$

- The optimal bandwidth choice is*

$$h^* = \frac{1}{n^{1/3}} \left(\frac{6}{\int f'(u)^2 du} \right)^{1/3}$$

- With this choice of h^* , we have*

$$\text{MISE}(\hat{f}) \approx \frac{C}{n^{2/3}} \quad \text{where} \quad C = \left(\frac{3}{4} \right)^{2/3} \left(\int f'(u)^2 du \right)^{1/3}.$$

Cross-validation

- Issue: the optimal bandwidth h^* depends on the unknown density f
- Idea: estimate the risk under each bandwidth selection h

$$L(h) := \int (\hat{f}(x) - f(x))^2 dx = \underbrace{\int \hat{f}^2(x) dx - 2 \int \hat{f}(x) f(x) dx}_{=: J(h)} + \int f^2(x) dx$$

- Cross-validation estimate of the risk:

$$\hat{J}(h) := \int \hat{f}^2(x) dx - \frac{2}{n} \sum_{i=1}^n \hat{f}_{(-i)}(X_i)$$

- It can be shown that $\hat{J}(h) \approx \mathbb{E}[J(h)]$
- **Cross validation:** select h that minimizes $\hat{J}(h)$

Cross-validation

- Issue: the optimal bandwidth h^* depends on the unknown density f
- Idea: estimate the risk under each bandwidth selection h

$$L(h) := \int (\hat{f}(x) - f(x))^2 dx = \underbrace{\int \hat{f}^2(x) dx - 2 \int \hat{f}(x) f(x) dx}_{=: J(h)} + \int f^2(x) dx$$

- Cross-validation estimate of the risk:

$$\hat{J}(h) := \int \hat{f}^2(x) dx - \frac{2}{n} \sum_{i=1}^n \hat{f}_{(-i)}(X_i)$$

- It can be shown that $\hat{J}(h) \approx \mathbb{E}[J(h)]$
- **Cross validation:** select h that minimizes $\hat{J}(h)$
- HW: prove the formula below that allows efficient computation of $\hat{J}(h)$:

$$\hat{J}(h) = \frac{2}{(n-1)h} - \frac{n+1}{n-1} \sum_{j=1}^m \hat{p}_j^2$$

Limitation of the histogram method

- Histograms are discontinuous (not a continuous density)
- The convergence rate $O(n^{-2/3})$ is not ideal
- Complicated in higher dimension (number of bins will be exponential in dimension)
- A better solution: [kernel density estimation](#)

Kernel Density Estimation (KDE)

Kernel density estimator (KDE):

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} K\left(\frac{x - x_i}{h}\right)$$

where $K(\cdot)$ is a kernel function and $h > 0$ is the bandwidth

- **Kernel function:** any function $K(x) \geq 0$ that satisfies

$$\int K(x)dx = 1, \quad \int xK(x)dx = 0, \quad \int x^2K(x)dx > 0$$

- Common kernel function:

- Gaussian Kernel: $K(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$
- Epanechnikov kernel: $K(x) = \frac{3}{4\sqrt{5}} \left(1 - \frac{x^2}{5}\right) \mathbb{1}\{|x| < \sqrt{5}\}$

Optimal bandwidth

Theorem 2.6

Under some regularity conditions, we have

$$R(f, \hat{f}_n) \approx \frac{h^4}{4} \left(\int x^2 K(x) dx \right)^2 \int (f''(x))^2 dx + \frac{1}{nh} \int K^2(x) dx.$$

The optimal bandwidth is

$$h^* = \frac{1}{n^{1/5}} \left(\int x^2 K(x) dx \right)^{-2/5} \left(\int K^2(x) dx \right)^{1/5} \left(\int (f''(x))^2 dx \right)^{-1/5}$$

With this choice of bandwidth,

$$R(f, \hat{f}_n) \asymp \frac{1}{n^{4/5}}.$$

Cross-validation

Cross-validation: estimate the risk under each bandwidth selection h

$$L(h) := \int (\hat{f}(x) - f(x))^2 dx = \underbrace{\int \hat{f}^2(x) dx - 2 \int \hat{f}(x) f(x) dx}_{=: J(h)} + \int f^2(x) dx$$

- Estimating $J(h)$:

$$\hat{J}(h) := \int \hat{f}^2(x) dx - \frac{2}{n} \sum_{i=1}^n \hat{f}_{(-i)}(X_i)$$

It can be shown that $\mathbb{E}[\hat{J}(h)] = \mathbb{E}[J(h)]$

- Cross validation: select h that minimizes $\hat{J}(h)$
- An efficient formula for approximately computing $\hat{J}(h)$:

$$\hat{J}(h) \approx \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \frac{1}{h} K^*\left(\frac{X_i - X_j}{h}\right) + \frac{2}{nh} K(0)$$

where $K^*(x) = \int K(x - y) K(y) dy - 2K(x)$

Theoretical guarantees for cross validation

Theorem 2.7 (Stone's Theorem)

Suppose that f is bounded. Let \hat{f}_h be the KDE with bandwidth h , and let \hat{h} be the bandwidth chosen by cross-validation. Then

$$\frac{\text{MISE}(\hat{f}_{\hat{h}})}{\inf_h \text{MISE}(\hat{f}_h)} \xrightarrow{P} 1$$

as $n \rightarrow \infty$.

- Stone's theorem provides theoretical justification for using cross-validation to select bandwidth for KDE.

Implications

Faster convergence rate: $\underbrace{O(n^{-4/5})}_{\text{KDE}}$ vs. $\underbrace{O(n^{-2/3})}_{\text{histogram}}$

Extension to higher dimension: consider estimating a density f in \mathbb{R}^d

- kernel function K : symmetric density (e.g., density of $\mathcal{N}(0, I_d)$)
- KDE: for a symmetric, PSD bandwidth matrix $\mathbf{H} \in \mathbb{R}^{d \times d}$

$$\hat{f}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \det(\mathbf{H})^{-1/2} K(\mathbf{H}^{-1/2}(\mathbf{x} - \mathbf{x}_i))$$

- bandwidth selection: Silverman's rule of thumb

$$H_{i,i} = \left(\frac{4}{n(d+2)} \right)^{2/(d+4)} \sigma_i \quad (1 \leq i \leq d), \quad H_{i,j} = 0 \quad (i \neq j).$$

where σ_i^2 is the variance of the i -th variable.

- suffers from **curse of dimensionality** (error exponential in d)

Kernel Density Classifier

Bayes optimal classifier:

$$\mathbb{P}(Y = k \mid X = x) = \frac{\mathbb{P}(X = x \mid Y = k) \mathbb{P}(Y = k)}{\sum_{k'=1}^K \mathbb{P}(X = x \mid Y = k') \mathbb{P}(Y = k')}$$

- Construct a KDE $\hat{f}_k(x)$ for the conditional density $\mathbb{P}(X = x \mid Y = k)$ using data $\{x_i : y_i = k\}$ for each class $k \in \{1, \dots, K\}$,
- Estimate class priors $\mathbb{P}(y = k)$ with empirical frequency $\hat{\pi}_k = n_k/n$
- Kernel density classifier: for any input x , return

$$\arg \max_{1 \leq k \leq K} \hat{\mathbb{P}}(Y = k \mid X = x) := \frac{\hat{\pi}_k \hat{f}_k(x)}{\sum_{k'=1}^K \hat{\pi}_{k'} \hat{f}_{k'}(x)}$$

Kernel Density Classifier

Bayes optimal classifier:

$$\mathbb{P}(Y = k \mid X = x) = \frac{\mathbb{P}(X = x \mid Y = k) \mathbb{P}(Y = k)}{\sum_{k'=1}^K \mathbb{P}(X = x \mid Y = k') \mathbb{P}(Y = k')}$$

- Construct a KDE $\hat{f}_k(x)$ for the conditional density $\mathbb{P}(X = x \mid Y = k)$ using data $\{x_i : y_i = k\}$ for each class $k \in \{1, \dots, K\}$,
- Estimate class priors $\mathbb{P}(y = k)$ with empirical frequency $\hat{\pi}_k = n_k/n$
- Kernel density classifier: for any input x , return

$$\arg \max_{1 \leq k \leq K} \hat{\mathbb{P}}(Y = k \mid X = x) := \frac{\hat{\pi}_k \hat{f}_k(x)}{\sum_{k'=1}^K \hat{\pi}_{k'} \hat{f}_{k'}(x)}$$

- Issue: **curse of dimensionality**

Naive Bayes Classifier

- The Naive Bayes model assumes that given a class $Y = k$, the features X_1, \dots, X_d are conditionally independent.
- The class-conditional density $f_k(x) \equiv \mathbb{P}(X = x \mid Y = k)$ is given by:

$$f_k(x) = \prod_{j=1}^d f_{k,j}(x_j) \quad \text{where} \quad x = (x_1, \dots, x_d)$$

where $f_{k,j}(X_k)$ is the marginal density of X_j conditional on $Y = k$

- Naive Bayes classifier: for any input x , return

$$\arg \max_{1 \leq k \leq K} \hat{\mathbb{P}}(Y = k \mid X = x) := \frac{\hat{\pi}_k \hat{f}_k(x)}{\sum_{k'=1}^K \hat{\pi}_{k'} \hat{f}_{k'}(x)}$$

where $\hat{f}_k(x) = \prod_{j=1}^d \hat{f}_{k,j}(x_j)$.

- The estimate $\hat{f}_{k,j}$ for class-conditional marginal densities $f_{k,j}$ can be computed using e.g., one-dimensional KDE or histogram

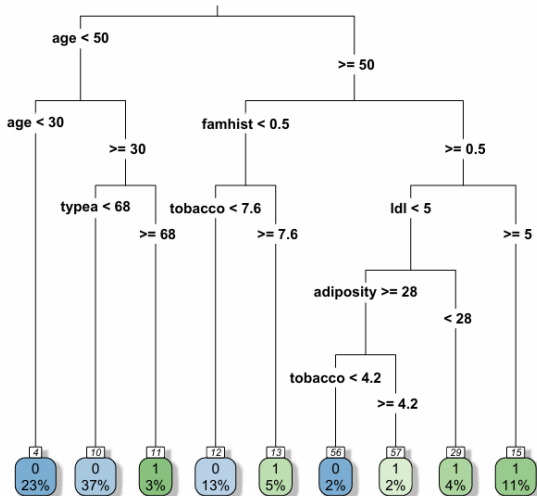
Discussions

- Naive Bayes works well in high-dimensional spaces and with small datasets, despite the independence assumption often being violated.
- **Advantages:**
 - Simple and fast
 - Avoids curse of dimensionality
 - Robust to irrelevant features
- **Disadvantages:**
 - Assumption of feature independence might be unrealistic

Tree-based methods

Classification tree

South African heart disease data: "0"="Yes, Disease", "1"="No"



Classification tree

Setup: $\mathcal{X} = \mathbb{R}^d$, $\mathcal{Y} = \{1, \dots, K\}$, training data $(X_1, Y_1), \dots, (X_n, Y_n)$

Idea: grow a tree to recursively partition the feature space into a set of rectangles, and do a simple majority vote in each rectangle

- Each node represents a rectangle in the feature space. The root node is the feature space $\mathcal{X} = \mathbb{R}^d$
- Each node is either a leaf (no children) or a parent (has two children)
- The left and right children comes from a partition of their parent node
- Suppose we have a collection of final partitioned regions associated with the leaves at the bottom of the tree, denoted by R_1, \dots, R_M
- For any input x , suppose that $x \in R_j$, then this classification tree returns

$$\hat{f}(x) = \arg \max_{k \in \mathcal{Y}} \sum_{X_i \in R_j} \mathbb{1}\{Y_i = k\}$$

i.e., the predicted label is the majority in the region R_j

How to grow a classification tree?

In order to grow a classification tree, we need to ask:

1. How to split each parent node?
2. How large should we grow the tree?

For the first question: **minimizing impurity**

- Suppose that the parent node is associated with a rectangle R
- Choose a covariate X_j and a split point t that minimizes the impurity
- Let the rectangles associated with its left and right children be

$$R_1(j, t) = \{X \in R : X_j \leq t\} \quad \text{and} \quad R_2(j, t) = \{X \in R : X_j > t\},$$

For the second question: **set some stopping criteria.**

- For example, we may fix some number n_0 , and we might stop partition a node when its associated rectangle has fewer than n_0 training data points.

Impurity function

Let R be the node to be split into two regions. We choose

$$\arg \min_{j,t} \underbrace{\frac{|R_1(j,t)|}{|R|} \gamma(R_1(j,t)) + \frac{|R_2(j,t)|}{|R|} \gamma(R_2(j,t))}_{\text{impurity function}},$$

- Here $\gamma(R)$ measures the “variance” of the labels of data in R : we want

$$\{Y_i : X_i \in R\} \quad \text{to have low variability}$$

- For any given rectangle R , let

$$p_k = \frac{1}{|R|} \sum_{X_i \in R} \mathbb{1}\{Y_i = k\}, \quad 1 \leq k \leq K.$$

Two common choice of the function $\gamma(\cdot)$:

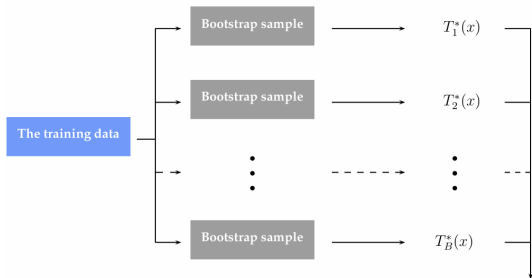
- **Gini index:** $\gamma(R) = \sum_k p_k(1 - p_k)$
- **Cross entropy:** $\gamma(R) = - \sum_k p_k \log p_k$

Insights

- **advantage:** the tree structure provides great **interpretability**
 - for example, it allows reasoning about the cause of diseases
- **disadvantage:** **instability** due to the use of *greedy* search:
 - splitting process is **greedy**
 - small changes in the training data can lead to **significantly different tree structures**
- **Solutions:**
 - Regularization: controlling tree growth parameters
 - Pruning: removing branches that do not provide significant predictive power
 - **Ensemble Methods:** use bagging to create a random forest

Bootstrap aggregating (Bagging)

- Training data $Z_n = \{(X_i, Y_i), 1 \leq i \leq n\}$
- Bootstrap sample $Z^{(*b)} = \{(X_i^{(*b)}, Y_i^{(*b)}), 1 \leq i \leq n\}$: sample n data points randomly from Z_n with replacement
- Apply the learning algorithm to the bootstrap sample for B times, and produce outcomes \hat{f}_b
- Majority vote: $\hat{f}^{\text{bagging}}(x) = \arg \max_{k \in \mathcal{Y}} \sum_{b=1}^B \mathbb{1}\{\hat{f}_b(x) = k\}$



Insights

- Trees generated in bagging are identically distributed (**not independent!**)
- Bias of bagged trees is the same as the individual tree
- **Pro:** Reduce the variance, so good for high-variance, low-bias procedures, like trees.
- **Heuristics:** Suppose we have B identically distributed random variables with variance σ^2 and positive pairwise correlation ρ , then their average has variance of

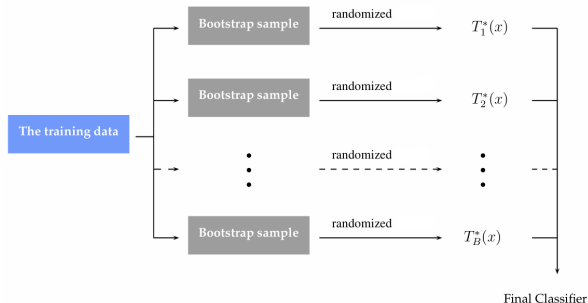
$$\rho\sigma^2 + \frac{1-\rho}{B}\sigma^2$$

- Increasing B does not reduce the first term

— Random Forest!

Random forests

- **Key idea:** use random dropout to decorrelate bootstrapped trees
- When growing a tree on a bootstrapped sample, before each split of the node, select $m \ll d$ variables at random as candidates to split
- Typical values for m is \sqrt{d} .
- Majority vote: $\hat{f}^{\text{RF}}(x) = \arg \max_{k \in \mathcal{Y}} \sum_{b=1}^B \mathbb{1}\{\tilde{f}_b(x) = k\}$



How to remove bias: Boosting

- Setup: $\mathcal{X} = \mathbb{R}^d$, $\mathcal{Y} = \{\pm 1\}$
- Weak classifier: error rate only slightly better than random guess
- Key idea: sequentially apply weak classification algorithm to repeatedly modified versions of the data to produce a sequence of weak classifiers
 - assign unequal weights to training data points — *possible for trees*
 - sequentially find a committee of weak classifiers $\{\hat{f}_m\}_{m=1}^M$
 - produce the final prediction through a weighted majority vote

$$\hat{f}(x) := \text{sign}\left(\sum_{m=1}^M \alpha_m \hat{f}_m(x)\right)$$

AdaBoost

Initialization: set the weights $w_i = 1/n$ for $1 \leq i \leq n$.

For $m = 1, \dots, M$:

- Fit a weak classifier $\hat{f}_m(x)$ using training data with weights $\omega_1, \dots, \omega_n$
- Compute the weighted misclassification error:

$$\text{err}^{(m)} = \frac{\sum_{i=1}^n w_i \mathbb{1}\{Y_i \neq \hat{f}_m(X_i)\}}{\sum_{i=1}^n w_i}.$$

- Compute:

$$\alpha_m = \log \left(\frac{1 - \text{err}^{(m)}}{\text{err}^{(m)}} \right).$$

- Update the weights by:

$$w_i \leftarrow w_i \cdot \exp \left(\alpha_m \cdot \mathbb{1}\{Y_i \neq \hat{f}_m(X_i)\} \right), \quad i = 1, 2, \dots, n.$$

Output: $\hat{f}(x) = \text{sign} \left(\sum_{m=1}^M \alpha_m \hat{f}_m(x) \right).$

AdaBoost: insights

Key idea: in the weight update step

$$w_i \leftarrow w_i \cdot \exp(\alpha_m \cdot \mathbb{1}\{Y_i \neq \hat{f}_m(X_i)\}), \quad i = 1, 2, \dots, n.$$

- For incorrectly classified data points, their weights get inflated by e^{α_m}
- Note that $\alpha_m > 0$ should always hold
- This re-weighting encourages the next classifier to focus more on the misclassified data points

Discussion: three main approaches to classification

Three main approaches

Bayes optimal classifier: for any $x \in \mathcal{X}$, output

$$f^*(x) := \arg \max_{y \in \mathcal{Y}} \mathbb{P}(Y = y \mid X = x)$$

minimizes the Bayes risk $R(f) = \mathbb{P}(f(X) \neq Y)$

Three main approaches

Bayes optimal classifier: for any $x \in \mathcal{X}$, output

$$f^*(x) := \arg \max_{y \in \mathcal{Y}} \mathbb{P}(Y = y \mid X = x)$$

minimizes the Bayes risk $R(f) = \mathbb{P}(f(X) \neq Y)$

- Plug-in approach: model data distribution ρ , then estimate densities

$$\mathbb{P}(Y = y \mid X = x) = \frac{\mathbb{P}(X = x \mid Y = y) \mathbb{P}(Y = y)}{\sum_{y' \in \mathcal{Y}} \mathbb{P}(X = x \mid Y = y') \mathbb{P}(Y = y')}$$

Three main approaches

Bayes optimal classifier: for any $x \in \mathcal{X}$, output

$$f^*(x) := \arg \max_{y \in \mathcal{Y}} \mathbb{P}(Y = y \mid X = x)$$

minimizes the Bayes risk $R(f) = \mathbb{P}(f(X) \neq Y)$

- Plug-in approach: model data distribution ρ , then estimate densities

$$\mathbb{P}(Y = y \mid X = x) = \frac{\hat{\mathbb{P}}(X = x \mid Y = y) \hat{\mathbb{P}}(Y = y)}{\sum_{y' \in \mathcal{Y}} \hat{\mathbb{P}}(X = x \mid Y = y') \hat{\mathbb{P}}(Y = y')}$$

Example: LDA, QDA, Kernel density classifier

Three main approaches

Bayes optimal classifier: for any $x \in \mathcal{X}$, output

$$f^*(x) := \arg \max_{y \in \mathcal{Y}} \mathbb{P}(Y = y \mid X = x)$$

minimizes the Bayes risk $R(f) = \mathbb{P}(f(X) \neq Y)$

- Plug-in approach
- Regression: modeling and estimating each

$$r_k(x) := \mathbb{P}(Y = k \mid X = x) \quad \text{for } k = 1, \dots, K$$

Three main approaches

Bayes optimal classifier: for any $x \in \mathcal{X}$, output

$$f^*(x) := \arg \max_{y \in \mathcal{Y}} \mathbb{P}(Y = y \mid X = x)$$

minimizes the Bayes risk $R(f) = \mathbb{P}(f(X) \neq Y)$

- Plug-in approach
- Regression: modeling and estimating each

$$r_k(x) := \mathbb{P}(Y = k \mid X = x) \quad \text{for } k = 1, \dots, K$$

Example: logistic regression

Three main approaches

Bayes optimal classifier: for any $x \in \mathcal{X}$, output

$$f^*(x) := \arg \max_{y \in \mathcal{Y}} \mathbb{P}(Y = y \mid X = x)$$

minimizes the Bayes risk $R(f) = \mathbb{P}(f(X) \neq Y)$

- Plug-in approach
- Regression
- Empirical risk minimization: choose a set of classifiers \mathcal{F} and find $\hat{f} \in \mathcal{F}$ that minimizes the “empirical risk”:

$$R_n(f) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{f(X_i) \neq Y_i\}$$

Intuition: when n is large, $R_n(f) \approx R(f)$ by LLN

Three main approaches

Bayes optimal classifier: for any $x \in \mathcal{X}$, output

$$f^*(x) := \arg \max_{y \in \mathcal{Y}} \mathbb{P}(Y = y \mid X = x)$$

minimizes the Bayes risk $R(f) = \mathbb{P}(f(X) \neq Y)$

- Plug-in approach
- Regression
- Empirical risk minimization
- Other approaches: SVM, tree-based methods...

ERM: advantages

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{f(X_i) \neq Y_i\} =: R_n(f)$$

- a straightforward method based on simple heuristics
- can be easily generalized to other loss $\ell(\cdot, \cdot)$ by considering

— *robustness!*

$$R_n(f) := \frac{1}{n} \sum_{i=1}^n \ell(f(X_i), Y_i)$$

if the ultimate goal is to minimize $R_\ell(f) = \mathbb{E}[\ell(f(X), Y)]$. For example, in binary classification (i.e., $\mathcal{Y} = \{0, 1\}$)

- Hinge loss $\ell(f(x), y) = \max\{0, 1 - yf(x)\}$
- Logistic loss $\ell(f(x), y) = \log(1 + \exp(-yf(x)))$

— *Logistic regression can also be viewed as ERM!*

ERM: disadvantages

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{f(X_i) \neq Y_i\} =: R_n(f)$$

- Not easy to compute (due to nonsmoothness of the indicator function)
- Solution: in binary classification (i.e., $\mathcal{Y} = \{0, 1\}$), consider using hinge loss or logistic loss $\ell(\cdot)$

$$R_n(f) := \frac{1}{n} \sum_{i=1}^n \ell(f(X_i), Y_i)$$

and relax $f : \mathbb{R}^d \rightarrow \mathbb{R}$, and finally output $\text{sign}(2(f(x) - 1))$

- Here we will only focus on the standard ERM

ERM: error decomposition

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{f(X_i) \neq Y_i\} =: R_n(f)$$

- We want to control the excess risk

$$R(\hat{f}_n) - R(f^*) = \underbrace{R(\hat{f}_n) - \min_{f \in \mathcal{F}} R(f)}_{\geq 0, \text{ statistical error}} + \underbrace{\min_{f \in \mathcal{F}} R(f) - R(f^*)}_{\geq 0, \text{ approximation error}}$$

- **approximation error**: becomes smaller when choosing larger \mathcal{F}
— *becomes 0 when $f^* \in \mathcal{F}$*
- **statistical error**: becomes smaller when n becomes larger, and when choosing smaller \mathcal{F} (why?)
- **trade-off between fit and complexity**
- In this course, we will focus on understanding **statistical error** with a given \mathcal{F} that includes f^* (so that **approximation error** = 0)

Excess risk via uniform deviations

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{f(X_i) \neq Y_i\} =: R_n(f)$$

Theorem 2.8

The excess risk is upper bounded by

$$R(\hat{f}_n) - R(f^*) \leq 2 \sup_{f \in \mathcal{F}} |R_n(f) - R(f)|$$

Implications:

- For a given f , we know that $R_n(f) \rightarrow R(f)$ at a rate $O(1/\sqrt{n})$ by CLT

$$\sqrt{n}(R_n(f) - R(f)) \xrightarrow{d} \mathcal{N}(0, \text{var}(\mathbb{1}\{f(X) \neq Y\}))$$

- But what about the **uniform convergence** of $\sup_{f \in \mathcal{F}} |R_n(f) - R(f)|$?

Concentration inequalities and uniform convergence

Why concentration inequalities?

Consider i.i.d. variables X_1, \dots, X_n with $\mathbb{E}[X_i] = \mu$ and $\text{var}(X_i) = \sigma^2$

- Central limit theorem (CLT):

$$\sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n X_i - \mu \right) \xrightarrow{d} \mathcal{N}(0, \sigma^2)$$

tells us that the sample average concentrates around μ , and the deviation scales like σ/\sqrt{n} as $n \rightarrow \infty$

- But this does not say anything useful when n is finite
- We want some non-asymptotic statement like:

$$\mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n X_i - \mu \right| \geq \varepsilon(n, \delta) \right) \leq \delta$$

holds for any $\delta > 0$, where $\varepsilon(n, \delta) > 0$ is some quantity that depends on the sample size n and the exceptional probability δ

A simple case with i.i.d. Gaussian

Suppose that $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu, \sigma^2)$, then we have

$$\frac{1}{n} \sum_{i=1}^n X_i - \mu \sim \mathcal{N}\left(0, \frac{\sigma^2}{n}\right)$$

Theorem 2.9

For $G \sim \mathcal{N}(0, 1)$ and any $t > 0$, we have

$$\left(\frac{1}{t} - \frac{1}{t^3}\right) \frac{1}{\sqrt{2\pi}} e^{-t^2/2} \leq \mathbb{P}(G \geq t) \leq \frac{1}{t} \frac{1}{\sqrt{2\pi}} e^{-t^2/2}.$$

As a result,

$$\mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - \mu\right| \geq t\right) \leq \frac{2\sigma}{\sqrt{nt}} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{nt^2}{2\sigma^2}\right)$$

A simple case with i.i.d. Gaussian

Suppose that $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu, \sigma^2)$, then we have

$$\frac{1}{n} \sum_{i=1}^n X_i - \mu \sim \mathcal{N}\left(0, \frac{\sigma^2}{n}\right)$$

Theorem 2.9

For $G \sim \mathcal{N}(0, 1)$ and any $t > 0$, we have

$$\left(\frac{1}{t} - \frac{1}{t^3}\right) \frac{1}{\sqrt{2\pi}} e^{-t^2/2} \leq \mathbb{P}(G \geq t) \leq \frac{1}{t} \frac{1}{\sqrt{2\pi}} e^{-t^2/2}.$$

As a result,

$$\mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - \mu\right| \geq t\right) \leq \frac{2\sigma}{\sqrt{nt}} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{nt^2}{2\sigma^2}\right)$$

Question: how can we extend these to more general distributions?

From Gaussian to sub-Gaussian

- Question: can we generalize these results to other random variables?
- Idea: consider other random variables with similar tail probability
- From Theorem 3.2, we know that for $G \sim \mathcal{N}(0, \sigma^2)$,

$$\mathbb{P}(|G| \geq t) \lesssim e^{-t^2/\sigma^2} \quad \text{for all } t \geq 0$$

- We may consider random variables satisfy this type of tail properties
— *sub-Gaussian*

Sub-Gaussian properties

Let X be a random variable, then the following properties are equivalent:

1. The tails of X satisfy

$$\mathbb{P}(|X| \geq t) \leq 2 \exp(-t^2/K_1^2) \quad \text{for all } t \geq 0$$

2. The moments of X satisfy

$$\|X\|_{L^p} := (\mathbb{E}[|X|^p])^{1/p} \leq K_2 \sqrt{p} \quad \text{for all } p \geq 1$$

3. The moment generating function (MGF) of X^2 satisfies

$$\mathbb{E}[\exp(\lambda^2 X^2)] \leq \exp(K_3^2 \lambda^2) \quad \text{for all } \lambda \text{ such that } |\lambda| \leq 1/K_3$$

4. The MGF of X^2 is bounded at some point, namely

$$\mathbb{E}[\exp(X^2/K_4^2)] \leq 2.$$

5. If $\mathbb{E}X = 0$, then the MGF of X satisfies

$$\mathbb{E}[\exp(\lambda X)] \leq \exp(K_5^2 \lambda^2) \quad \text{for all } \lambda \in \mathbb{R}.$$

where $K_1, \dots, K_5 > 0$ may differ by at most a multiplicative constant factor

Sub-Gaussian distributions: definition

- If X satisfies one of properties 1-4, it is a *sub-Gaussian random variable*.
- The *sub-Gaussian norm* of X , denoted $\|X\|_{\psi_2}$, is defined to be the smallest K_4 in property 4. In other words, we define

$$\|X\|_{\psi_2} = \inf \left\{ t > 0 : \mathbb{E} \exp \left(X^2 / t^2 \right) \leq 2 \right\}.$$

— can also be defined using K_1 , K_2 or K_3

- Properties: there exists some absolute constants $c, C > 0$ such that
 - $P(|X| \geq t) \leq 2 \exp \left(- ct^2 / \|X\|_{\psi_2}^2 \right)$
 - $\|X\|_{L^p} \leq C \|X\|_{\psi_2} \sqrt{p}$
 - $\mathbb{E} \exp \left(X^2 / \|X\|_{\psi_2}^2 \right) \leq 2$
 - if $\mathbb{E}[X] = 0$, then $\mathbb{E} \exp(\lambda X) \leq \exp(C \lambda^2 \|X\|_{\psi_2}^2)$

Sub-Gaussian distributions: examples

- **Gaussian:** if $X \sim \mathcal{N}(0, \sigma^2)$, then X is sub-Gaussian with

$$\|X\|_{\psi_2} \leq C\sigma$$

for some universal constant $C = 2\sqrt{2/3}$.

- **Bounded:** any bounded random variable X is sub-Gaussian with

$$\|X\|_{\psi_2} \leq C\|X\|_{\infty}$$

for some universal constant $C = 1/\sqrt{\log 2}$.

Sub-Gaussian norm can be viewed as a characterization of “magnitude” for light tail distributions.

Centering and independent sums

Theorem 2.10

- If X is sub-Gaussian, then $X - \mathbb{E}[X]$ is sub-Gaussian with

$$\|X - \mathbb{E}[X]\|_{\psi_2} \leq C\|X\|_{\psi_2}$$

where C is an absolute constant.

- Let X_1, \dots, X_N be independent, mean zero, sub-Gaussian random variables. Then the sum $S_N = \sum_{i=1}^N X_i$ is also sub-Gaussian, and its sub-Gaussian norm satisfies

$$\|S_N\|_{\psi_2}^2 \leq C \sum_{i=1}^N \|X_i\|_{\psi_2}^2,$$

where C is an absolute constant.

Analog:

- If X_1, \dots, X_n are i.i.d. $\mathcal{N}(0, \sigma^2)$, then $S_N \sim \mathcal{N}(0, N\sigma^2)$
- If X_1, \dots, X_n are independent with $\|X_i\|_{\psi_2} \leq \sigma$, then $\|S_N\|_{\psi_2} \lesssim \sqrt{N}\sigma$

Hoeffding's inequality

Theorem 2.11 (Hoeffding's Inequality)

Let X_1, \dots, X_N be independent, mean-zero, sub-Gaussian random variables. Then, for any $t \geq 0$, we have:

$$\mathbb{P} \left(\left| \sum_{i=1}^N X_i \right| \geq t \right) \leq 2 \exp \left(- \frac{ct^2}{\sum_{i=1}^N \|X_i\|_{\psi_2}^2} \right),$$

where c is an absolute constant.

Implications

- **General Hoeffding:** under the setup of Theorem 3.4, consider any vector $\mathbf{a} = (a_1, \dots, a_n) \in \mathbb{R}^n$, we have

$$\mathbb{P} \left(\left| \sum_{i=1}^N a_i X_i \right| \geq t \right) \leq 2 \exp \left(- \frac{ct^2}{K^2 \|\mathbf{a}\|_2^2} \right),$$

where $K := \max \|X_i\|_{\psi_2}$.

- **Example:** suppose that $X_i \sim \text{Bernoulli}(p_i)$ for $1 \leq i \leq n$, then

$$\mathbb{P} \left(\left| \sum_{i=1}^N (X_i - p_i) \right| \geq t \right) \leq 2 \exp \left(- \frac{ct^2}{N} \right),$$

A sharper result for binomial concentration: Chernoff's inequality (HW)

Back to ERM: finite \mathcal{F}

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{f(X_i) \neq Y_i\} =: R_n(f)$$

Theorem 2.12

Suppose that \mathcal{F} is a finite set. Then with probability exceeding $1 - \delta$, the excess risk of ERM is upper bounded by

$$R(\hat{f}_n) - R(f^*) \leq C \sqrt{\frac{\log(|\mathcal{F}|/\delta)}{n}}.$$

for some universal constant $C > 0$.

- Key proof idea: **union bound argument**
- What if \mathcal{F} is not finite (e.g., the set of linear classifiers)?

— use *VC dimension!*

Back to ERM: finite \mathcal{F}

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{f(X_i) \neq Y_i\} =: R_n(f)$$

Theorem 2.12

Suppose that \mathcal{F} is a finite set. Then with probability exceeding $1 - \delta$, the excess risk of ERM is upper bounded by

$$R(\hat{f}_n) - R(f^*) \leq C \sqrt{\frac{\log(|\mathcal{F}|/\delta)}{n}}.$$

for some universal constant $C > 0$.

- Key proof idea: **union bound argument**
- What if \mathcal{F} is not finite (e.g., the set of linear classifiers)?
— use *VC dimension!*
- But before going into that, let's first warm up with something simpler

ℓ_2 norm of a sub-Gaussian random vector

- Consider a random vector $\mathbf{x} = (X_1, \dots, X_d)$, where X_1, \dots, X_d are independent random variables with $\mathbb{E}[X_i] = 0$ and $\|X_i\|_{\psi_2} \leq \sigma$
- Can we establish a non-asymptotic upper bound for $\|\mathbf{x}\|_2$?

ℓ_2 norm of a sub-Gaussian random vector

- Consider a random vector $\mathbf{x} = (X_1, \dots, X_d)$, where X_1, \dots, X_d are independent random variables with $\mathbb{E}[X_i] = 0$ and $\|X_i\|_{\psi_2} \leq \sigma$
- Can we establish a non-asymptotic upper bound for $\|\mathbf{x}\|_2$?
- Solution 1: entrywise concentration and union bound

$$\mathbb{P}(\|\mathbf{x}\|_2 \leq C\sigma\sqrt{d\log(d/\delta)}) \geq 1 - \delta$$

for some universal constant $C > 0$

ℓ_2 norm of a sub-Gaussian random vector

- Consider a random vector $\mathbf{x} = (X_1, \dots, X_d)$, where X_1, \dots, X_d are independent random variables with $\mathbb{E}[X_i] = 0$ and $\|X_i\|_{\psi_2} \leq \sigma$
- Can we establish a non-asymptotic upper bound for $\|\mathbf{x}\|_2$?
- Solution 1: entrywise concentration and union bound

$$\mathbb{P}(\|\mathbf{x}\|_2 \leq C\sigma\sqrt{d\log(d/\delta)}) \geq 1 - \delta$$

for some universal constant $C > 0$

- Solution 2: uniform concentration using

$$\|\mathbf{x}\|_2 = \sup_{\mathbf{a} \in \mathcal{S}^{d-1}} \mathbf{a}^\top \mathbf{x}$$

where $\mathcal{S}^{d-1} := \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 = 1\}$ is the unit sphere in \mathbb{R}^d

— *could this provide a better concentration bound?*

Operator norm of sub-Gaussian matrix

- Consider a random matrix $\mathbf{X} = (X_{i,j})_{1 \leq i,j \leq d}$ with independent entries that satisfies $\mathbb{E}[X_{i,j}] = 0$ and $\|X_{i,j}\|_{\psi_2} \leq \sigma$
- Can we establish a non-asymptotic upper bound for $\|\mathbf{X}\|$?
- Operator norm:

$$\|\mathbf{X}\| = \sup_{\mathbf{a} \in \mathcal{S}^{d-1}} \|\mathbf{X}\mathbf{a}\|_2 = \sup_{\mathbf{a}, \mathbf{b} \in \mathcal{S}^{d-1}} \mathbf{a}^\top \mathbf{X} \mathbf{b}$$

where $\mathcal{S}^{d-1} := \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 = 1\}$ is the unit sphere in \mathbb{R}^d

A framework for uniform concentration

- **Goal:** upper bounding $\sup_{\mathbf{a} \in \mathcal{S}^{d-1}} \mathbf{a}^\top \mathbf{x}$
- **Step 1: pointwise concentration.** For any fixed $\mathbf{a} \in \mathcal{S}^{d-1}$, we can use Hoeffding's inequality to get

$$\mathbb{P}(|\mathbf{a}^\top \mathbf{x}| \leq C\sigma\sqrt{\log(1/\delta)}) \geq 1 - \delta$$

for some universal constant $C > 0$

- **Difficulty:** the unit sphere \mathcal{S}^{d-1} is not a finite set, union bound argument does not work
- **Idea:** find a finite subset \mathcal{N} of \mathcal{S}^{d-1} that is *fine* enough, such that

$$\sup_{\mathbf{a} \in \mathcal{S}^{d-1}} \mathbf{a}^\top \mathbf{x} \stackrel{?}{\lesssim} \sup_{\mathbf{a} \in \mathcal{N}} \mathbf{a}^\top \mathbf{x} \leq C\sigma\sqrt{\log(|\mathcal{N}|/\delta)}$$

with probability at least $1 - \delta$

Epsilon net

- Let (T, d) be a metric space. Consider a subset $K \subset T$ and let $\varepsilon > 0$.
 - *e.g., consider $T = \mathbb{R}^d$, $d(\cdot, \cdot)$ is Euclidean distance, $K = \mathcal{S}^{d-1}$*
- A subset $N \subseteq K$ is called an ε -net of K if every point in K is within distance ε of some point of N , i.e.,

$$\forall x \in K, \quad \exists x_0 \in N \quad \text{s.t.} \quad d(x, x_0) \leq \varepsilon.$$

Theorem 2.13

Let \mathcal{N}_ε be an ε -net of \mathcal{S}^{d-1} . If $\varepsilon < 1$, then for any $\mathbf{x} \in \mathbb{R}^d$,

$$\sup_{\mathbf{a} \in \mathcal{N}_\varepsilon} \mathbf{a}^\top \mathbf{x} \leq \sup_{\mathbf{a} \in \mathcal{S}^{d-1}} \mathbf{a}^\top \mathbf{x} \leq \frac{1}{1 - \varepsilon} \sup_{\mathbf{a} \in \mathcal{N}_\varepsilon} \mathbf{a}^\top \mathbf{x},$$

and if $\varepsilon < 1/2$, then for any $\mathbf{X} \in \mathbb{R}^{d \times d}$,

$$\sup_{\mathbf{a}, \mathbf{b} \in \mathcal{N}_\varepsilon} \mathbf{a}^\top \mathbf{X} \mathbf{b} \leq \sup_{\mathbf{a}, \mathbf{b} \in \mathcal{S}^{d-1}} \mathbf{a}^\top \mathbf{X} \mathbf{b} \leq \frac{1}{1 - 2\varepsilon} \sup_{\mathbf{a}, \mathbf{b} \in \mathcal{N}_\varepsilon} \mathbf{a}^\top \mathbf{X} \mathbf{b}.$$

The covering number

Covering number: the **smallest** possible cardinality of an ε -net of K , denoted by $\mathcal{N}(K, \varepsilon)$

Theorem 2.14

The covering number of \mathcal{S}^{d-1} is upper bounded by

$$\mathcal{N}(\mathcal{S}^{d-1}, \varepsilon) \leq \left(\frac{2}{\varepsilon} + 1 \right)^d$$

ℓ_2 norm of sub-Gaussian random vector

- **Goal:** upper bounding $\sup_{\mathbf{a} \in \mathcal{S}^{d-1}} \mathbf{a}^\top \mathbf{x}$
- **Step 1: pointwise concentration.** For any fixed $\mathbf{a} \in \mathcal{S}^{d-1}$, we can use Hoeffding's inequality to get

$$\mathbb{P}(|\mathbf{a}^\top \mathbf{x}| \leq C_1 \sigma \sqrt{\log(1/\delta)}) \geq 1 - \delta$$

for some universal constant $C_1 > 0$

ℓ_2 norm of sub-Gaussian random vector

- **Goal:** upper bounding $\sup_{\mathbf{a} \in \mathcal{S}^{d-1}} \mathbf{a}^\top \mathbf{x}$
- **Step 1: pointwise concentration.** For any fixed $\mathbf{a} \in \mathcal{S}^{d-1}$, we can use Hoeffding's inequality to get

$$\mathbb{P}(|\mathbf{a}^\top \mathbf{x}| \leq C_1 \sigma \sqrt{\log(1/\delta)}) \geq 1 - \delta$$

for some universal constant $C_1 > 0$

- **Step 2: uniform concentration over an $1/2$ -net.** Let $\mathcal{N}_{1/2}$ be the smallest $1/2$ -net of \mathcal{S}^{d-1} . By union bound argument and Theorem 3.7,

$$\mathbb{P}\left(\sup_{\mathbf{a} \in \mathcal{N}_{1/2}} |\mathbf{a}^\top \mathbf{x}| \leq C_2 \sigma \sqrt{d \log(1/\delta)}\right) \geq 1 - \delta$$

for some universal constant $C_2 > 0$

ℓ_2 norm of sub-Gaussian random vector

- **Goal:** upper bounding $\sup_{\mathbf{a} \in \mathcal{S}^{d-1}} \mathbf{a}^\top \mathbf{x}$
- **Step 1: pointwise concentration.** For any fixed $\mathbf{a} \in \mathcal{S}^{d-1}$, we can use Hoeffding's inequality to get

$$\mathbb{P}(|\mathbf{a}^\top \mathbf{x}| \leq C_1 \sigma \sqrt{\log(1/\delta)}) \geq 1 - \delta$$

for some universal constant $C_1 > 0$

- **Step 2: uniform concentration over an $1/2$ -net.** Let $\mathcal{N}_{1/2}$ be the smallest $1/2$ -net of \mathcal{S}^{d-1} . By union bound argument and Theorem 3.7,

$$\mathbb{P}\left(\sup_{\mathbf{a} \in \mathcal{N}_{1/2}} |\mathbf{a}^\top \mathbf{x}| \leq C_2 \sigma \sqrt{d \log(1/\delta)}\right) \geq 1 - \delta$$

for some universal constant $C_2 > 0$

- **Step 3: approximation.** By Theorem 3.6,

$$\mathbb{P}(\|\mathbf{x}\|_2 \leq C_3 \sigma \sqrt{d \log(1/\delta)}) \geq 1 - \delta$$

for some universal constant $C_3 > 0$

Operator norm of sub-Gaussian random matrix

- **Goal:** upper bounding $\sup_{\mathbf{a}, \mathbf{b} \in \mathcal{S}^{d-1}} \mathbf{a}^\top \mathbf{X} \mathbf{b}$
- **Step 1: pointwise concentration.** For any fixed $\mathbf{a}, \mathbf{b} \in \mathcal{S}^{d-1}$, we can use Hoeffding's inequality to get

$$\mathbb{P}(|\mathbf{a}^\top \mathbf{X} \mathbf{b}| \leq C_1 \sigma \sqrt{\log(1/\delta)}) \geq 1 - \delta$$

for some universal constant $C_1 > 0$

Operator norm of sub-Gaussian random matrix

- **Goal:** upper bounding $\sup_{\mathbf{a}, \mathbf{b} \in \mathcal{S}^{d-1}} \mathbf{a}^\top \mathbf{X} \mathbf{b}$
- **Step 1: pointwise concentration.** For any fixed $\mathbf{a}, \mathbf{b} \in \mathcal{S}^{d-1}$, we can use Hoeffding's inequality to get

$$\mathbb{P}(|\mathbf{a}^\top \mathbf{X} \mathbf{b}| \leq C_1 \sigma \sqrt{\log(1/\delta)}) \geq 1 - \delta$$

for some universal constant $C_1 > 0$

- **Step 2: uniform concentration over an $1/4$ -net.** Let $\mathcal{N}_{1/4}$ be the smallest $1/4$ -net of \mathcal{S}^{d-1} . By union bound argument and Theorem 3.7,

$$\mathbb{P}\left(\sup_{\mathbf{a}, \mathbf{b} \in \mathcal{N}_{1/4}} |\mathbf{a}^\top \mathbf{X} \mathbf{b}| \leq C_2 \sigma \sqrt{d \log(1/\delta)}\right) \geq 1 - \delta$$

for some universal constant $C_2 > 0$

Operator norm of sub-Gaussian random matrix

- **Goal:** upper bounding $\sup_{\mathbf{a}, \mathbf{b} \in \mathcal{S}^{d-1}} \mathbf{a}^\top \mathbf{X} \mathbf{b}$
- **Step 1: pointwise concentration.** For any fixed $\mathbf{a}, \mathbf{b} \in \mathcal{S}^{d-1}$, we can use Hoeffding's inequality to get

$$\mathbb{P}(|\mathbf{a}^\top \mathbf{X} \mathbf{b}| \leq C_1 \sigma \sqrt{\log(1/\delta)}) \geq 1 - \delta$$

for some universal constant $C_1 > 0$

- **Step 2: uniform concentration over an $1/4$ -net.** Let $\mathcal{N}_{1/4}$ be the smallest $1/4$ -net of \mathcal{S}^{d-1} . By union bound argument and Theorem 3.7,

$$\mathbb{P}\left(\sup_{\mathbf{a}, \mathbf{b} \in \mathcal{N}_{1/4}} |\mathbf{a}^\top \mathbf{X} \mathbf{b}| \leq C_2 \sigma \sqrt{d \log(1/\delta)}\right) \geq 1 - \delta$$

for some universal constant $C_2 > 0$

- **Step 3: approximation.** By Theorem 3.6,

$$\mathbb{P}(\|\mathbf{X}\|_2 \leq 2C_2 \sigma \sqrt{d \log(1/\delta)}) \geq 1 - \delta$$

for some universal constant $C_3 > 0$

VC dimension

- Let \mathcal{F} be a class of binary functions on the domain \mathcal{X} .
- **Shattering:** a set of points $\{x_1, \dots, x_k\} \subseteq \mathcal{X}$ is shattered by \mathcal{F} if for every possible labeling $\{0, 1\}^k$, there exists a function $f \in \mathcal{F}$ that realizes the labeling.
- The **VC dimension** of \mathcal{F} , denoted $\text{VC}(\mathcal{F})$, is the largest integer k such that there exists a set of k points in \mathcal{X} that can be *shattered* by \mathcal{F} .
- Examples:
 - When $\mathcal{X} = \mathbb{R}^2$, \mathcal{F} = linear classifiers, we have $\text{vc}(\mathcal{F}) = 3$
 - In general, when $\mathcal{X} = \mathbb{R}^d$, \mathcal{F} = linear classifiers, then $\text{vc}(\mathcal{F}) = d + 1$

Bounding excess risk via VC dimension

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{f(X_i) \neq Y_i\} =: R_n(f)$$

Theorem 2.15

Suppose that \mathcal{F} is a class of Boolean function with $\text{vc}(\mathcal{F}) < \infty$. Then with probability exceeding $1 - \delta$,

$$R(\hat{f}_n) - R(f^*) \leq C \sqrt{\frac{\text{vc}(\mathcal{F}) \log(1/\delta)}{n}}$$

for some universal constant $C > 0$.

Implications:

- For $\mathcal{F} =$ linear classifiers in \mathbb{R}^d , the excess risk is $O(\sqrt{d/n})$.