# Lecture 6

Classifiers:
- Generative vs discriminative
- Loss functions
    empirical loss, Bayes loss

- Bias & Variance — for parameter estimation
                    — what's different for prediction

Q1  Thursday
    beginning of class

HW2 — due Monday
    10/23   11:59pm

# Lecture II: Prediction – Basic concepts

Marina Meilă
mmp@stat.washington.edu

Department of Statistics
University of Washington

October, 2023

Parametric vs non-parametric

**Generative and discriminative models for classification**
  Generative classifiers
  Discriminative classifiers
  Generative vs discriminative classifiers

Loss functions
  Bayes loss

Variance, bias and complexity

**Reading** HTF Ch.: 2.1–5,2.9, 7.1–4 bias-variance tradeoff, Murphy Ch.: 1., 8.6[1], Bach Ch.:

---

[1]Neither textbook is close to these notes except in a few places; take them as alternative perspectives or related reading

# The "learning" problem

▶ **Given**
▶ a problem (e.g. recognize digits from $m \times m$ gray-scale images)
▶ a **sample** or (**training set**) of **labeled data**

$$\mathcal{D} = \{(x^1, y^1), (x^2, y^2), \ldots (x^n, y^n)\}$$

drawn i.i.d. from an unknown $P_{XY}$
▶ **model class** $\mathcal{F} = \{f\}$ = set of predictors to choose from

▶ **Wanted**
▶ a predictor $f \in \mathcal{F}$ that performs well on future samples from the same $P_{XY}$

▶ "choose a predictor $f \in \mathcal{F}$" = training/learning
▶ "performs well on future samples" (i.e. $f$ **generalizes** well) – how do we measure this? how can we "guarantee" it?
▶ choosing $\mathcal{F}$ is the **model selection problem** – about this later

# A zoo of predictors

- Linear regression
- Logistic regression
- Linear Discriminant (LDA) ✓
- Quadratic Discriminant (QDA) ✓
- CART (Decision Trees) **Discriminative**
- K-Nearest Neighbors ― " ―
- Nadaraya-Watson (Kernel regression) ← **NW classifier ?**
- Naive Bayes ✓ **generative**
- Neural networks/Deep learning **D**
- Support Vector Machines **D**
- Monotonic Regression

# Generative classifiers

One way to define a classifier is to assume that each class is generated by a distribution $g_y(X) = P(X|Y = y)$. If we know the distributions $g_y$ and the class probabilities $P(Y = y)$, we can derive the *posterior probability* distribution of $Y$ for a given $x$. This is

$$P(Y = y|X) = \frac{P(Y = y)g_y(X)}{\sum_{y'} P(Y = y')g_{y'}(X)} = \frac{P(Y = y)g_y(X)}{P(X)} \tag{1}$$

The "best guess" for $Y(X)$ (i.e. the decision rule) is

$$f(X) = \mathrm{argmax}_y P(Y = y|x) = \mathrm{argmax}_y P(Y = y)g_y(x) \tag{2}$$

► (1) amounts to a likelihood ratio test for $Y$.
► The functions $g_y(x)$ are known as **generative models** for the classes $y$.
   Therefore, the resulting classifier is called a **generative classifier**.
   Examples: LDA, QDA, Naive Bayes.
► In contrast, a classifier defined directly in terms of $f(x)$ (or $P_{Y|X}$), like the linear, quadratic, decision tree is called a **discriminative classifier**.
► In practice, we may not know the functions $g_y(x)$, in which case we estimate them from the sample $\mathcal{D}$.

# Generative classifiers (Binary)

given $\mathcal{D} = \{(x^i, y^i), i=1:n\}$    $y^i \in \{\pm 1\}$

1. fit   $g_+ \overset{\text{"="}}{=} P_{X|y=+}$    ← generative models

      $g_- \overset{\text{"="}}{=} P_{X|y=-}$

      ↑ density or PMF

Ex: Classify objects by images    $g_{car} = f(\text{image} \mid car)$    ← density estimation

        $g_{person} = f(-\text{"}- \mid person)$

2. $P_\pm = \dfrac{n_\pm}{n}$    $n_\pm = |\{x^i \mid y^i = \pm\}|$

             counts

3. Bayes' rule

        prior ←   ← likelihood

$$P[y=+ \mid x] = \frac{P_+ \, g_+(x)}{P_+ g_+(x) + P_- g_-(x)}$$

$\Rightarrow \hat{y}(x) = \underset{y}{\arg\max} \; \underline{P[y|x]}$

                              confidence

Ex: • Extends to any $y \in \{1, 2, \dots m\}$

              multiclass

# Generative classifier and the likelihood ratio

$$P(Y = y | X) = \frac{P(Y = y) g_y(X)}{\sum_{y'} P(Y = y') g_{y'}(X)} = \frac{P(Y = y) g_y(X)}{P(X)}$$

$$f(x) = \text{argmax}_y P(Y = y | x) = \text{argmax}_y g_y(x) P(Y = y)$$

Likelihood Ratio test (for $y \in \{\pm 1\}$)

$$\frac{g_+(x) P(Y = +)}{g_-(x) P(Y = -)} = \frac{P[+ | x]}{P[- | x]}$$

## Example (Fisher's LDA in one dimension)

Assume $Y = \pm 1$, $g_y(x) = N(x, \pm\mu, \sigma^2 I)$, i.e each class is generated by a Normal distribution with the same spherical covariance matrix, but with a different mean. Let $P(Y = 1) = p \in (0, 1)$. Then, the posterior probability of $Y$ is

$$P(Y = 1|x) \propto p e^{-||x-\mu||^2/(2\sigma^2)} \quad P(Y = -1|x) \propto (1-p)e^{-||x+\mu||^2/(2\sigma^2)} \quad (3)$$

and $f(x) = 1$ iff $\ln P(Y = 1|x)/P(Y = -1|x) \geq 0$, i.e iff

$$\ln \frac{p}{1-p} - \frac{1}{2\sigma^2}[||x^2|| - 2\mu^T x + ||\mu||^2 - ||x^2|| - (2\mu)^T x - ||\mu||^2] = \left(\frac{2\mu}{\sigma^2}\right)^T x + \ln \frac{p}{1-p} \geq 0 \quad (4)$$

Hence, the classifier $f(x)$ turns out to be a linear classifier. The decision boundary is perpendicular to the segment connecting the centers $\mu, -\mu$. This classifier is known as **Fisher's Linear Discriminant**. [**Exercises** Show that if the generative models are normal with different variances, then we obtain a quadratic classifier. What happens if the models $g_y$ have the same variance, but it is a full covariance matrix $\Sigma$?]

$$g_{\pm} = N(\mu_{\pm}, \sigma^2 I)$$

↳ same covariance

# Discriminative classifiers

▶ Defined directly in terms of $f(x)$ or (almost) equivalently, in terms of the decision boundary $\{f(x) = 0\}$
▶ Can be classified by the shape of the decision boundary (if it's simple)
  ▶ linear, polygonal, quadratic, cubic,...

## The ambiguity of "linear classifier"

Does it mean $f(x) = \beta^T x$ OR $\{f(x) = 0\}$ is a hyperplane ?
If we talk about **classification** and the domain of $x$ is $\mathbb{R}^d$, then "linear" refers to decision boundary. Otherwise it refers to the expression of $f(x)$. Exercise Find examples when the two definitions are not equivalent

▶ Can be grouped by model class (obviously)
  ▶ Neural network, K-nearest neighbor, decision tree, ...
    Exercise Is logistic regression a generative or discriminative classifier?
▶ By method of training (together with model class)
  ▶ For example, PERCEPTRON algorithm, Logistic Regression, (Linear) Support Vector Machine (see later), Decision Tree with 1 level are all linear classifiers, but usually produce different decision boundaries give a $\mathcal{D}$

# A comparison of generative and discriminative classifiers

### Advantages of generative classifiers
- Generative classifiers are statistically motivated ✔
- Generative classifiers are *asymptotically optimal* ✔

$$g_y = P(x \mid y)$$
$$g_y \longrightarrow \text{True } P_{X|Y}$$

## Theorem

*If $Y \in \{\pm 1\}$, the model class $G_y$ in which we are estimating $g_y$ contains the true distributions $P(X|Y = y)$ for every $y$, and $g_y = P(X|Y), P(Y = y)$ are estimated by Maximum Likelihood then the expected loss[2] of the generative classifier $f_g$ given by (2) tends to the Bayes loss when $n \to \infty$, i.e $\lim_{n\to\infty} L_{01}(f_g) \leq \min_{f \in \mathcal{F}} L_{01}(f)$. Here $\mathcal{F}$ is the class of likelihood ratio classifiers obtainable from $g_y$'s in $\mathcal{G}_y$.*
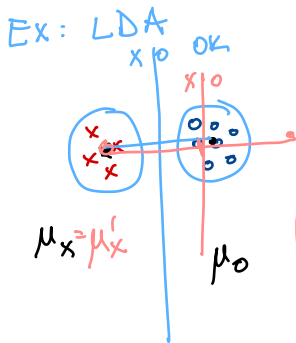
- The log-likelihood ratio $\ln \frac{P(Y=1|x)}{P(Y=-1|x)}$ is a natural confidence measure for the label at $f_g(x)$. The further away from 0 the likelihood ratio, the higher the confidence that the chosen $y$ is correct.
- Generative classifiers extend naturally to more than two classes. If a new class appears, or the class distribution $P(Y)$ changes, updating the classifier is simple and computationally efficient.
- Often it is easier to pick a (parametric) model class for $g_y$ than an $f$ directly. Generative models are generally more intuitive, while often representing/visualizing decision boundaries between more than two classes is tedious.

---

[2] Loss, Bayes loss, $L01$ are defined in the next section.

## Advantages of discriminative classifiers

▶ Generative models offer no guarantees if the true $g_y$ aren't in the chosen model class, whereas for many classes of discriminative models there are guarantees.

▶ Many discriminative models have performance guarantees for any sample size $n$, while generative models are only guaranteed for large enough $n$

▶ Discriminative classifiers offer many more choices (but one must know how to pick the right model)

▶ Generative models do not use data optimally in the non-asymptotic regime (when $n \ll \infty$). This has been confirmed practically many times, as discriminative classifiers have been very successful for limited sample sizes

Exercise LDA vs Logistic regression: Experiment with LDA vs LR when data comes from 2 Normal distributions, with outliers. What outliers affect which method more? Experiment also on a toy data set like the one in the lecture notes.



$$\text{Ex : LDA}$$

$$\text{LR}$$

$$\beta^T x = f(x) = \ln \frac{P[1|x]}{P[-1|x]}$$

models this

not these or $P_{x|y}$ separately

$\Rightarrow$ LR not generative

$\mu_x = \mu_x'$

$\mu_0$

$\mu_0'$

add outlier

# Loss functions

The **loss function** represents the cost of error in a prediction problem. We denote it by $L$, where

$L(y, \hat{y}) =$ the cost of predicting $\hat{y}$ when the actual outcome is $y$

*true* → $y$

↘ *predicted*

Note that sometimes the loss depends on $x$ directly. Then we would write it as $L(y, \hat{y}, x)$.

As usually $\hat{y} = f(x)$ or $\mathrm{sgn} f(x)$, we will typically abuse notation and write $L(y, f(x))$.

We assume $L \geq 0$ always

# Least Squares (LS) loss

## Regression

The **Least Squares (LS)** (or **quadratic**) loss function is given by

$$L_{LS}(y, f(x)) = \frac{1}{2}(y - f(x))^2 \qquad (5)$$

This loss is commonly associated with regression problems.

Example: $L_{LS}$ is the log-likelihood of a regression problem (linear or not) with Gaussian noise.
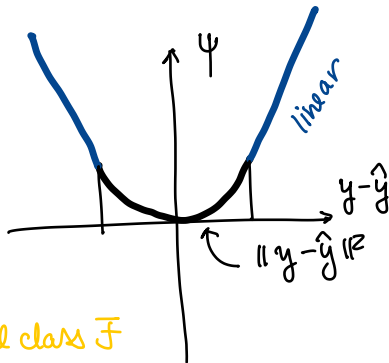
- $L_4(y, \hat{y}) = |y - \hat{y}|$

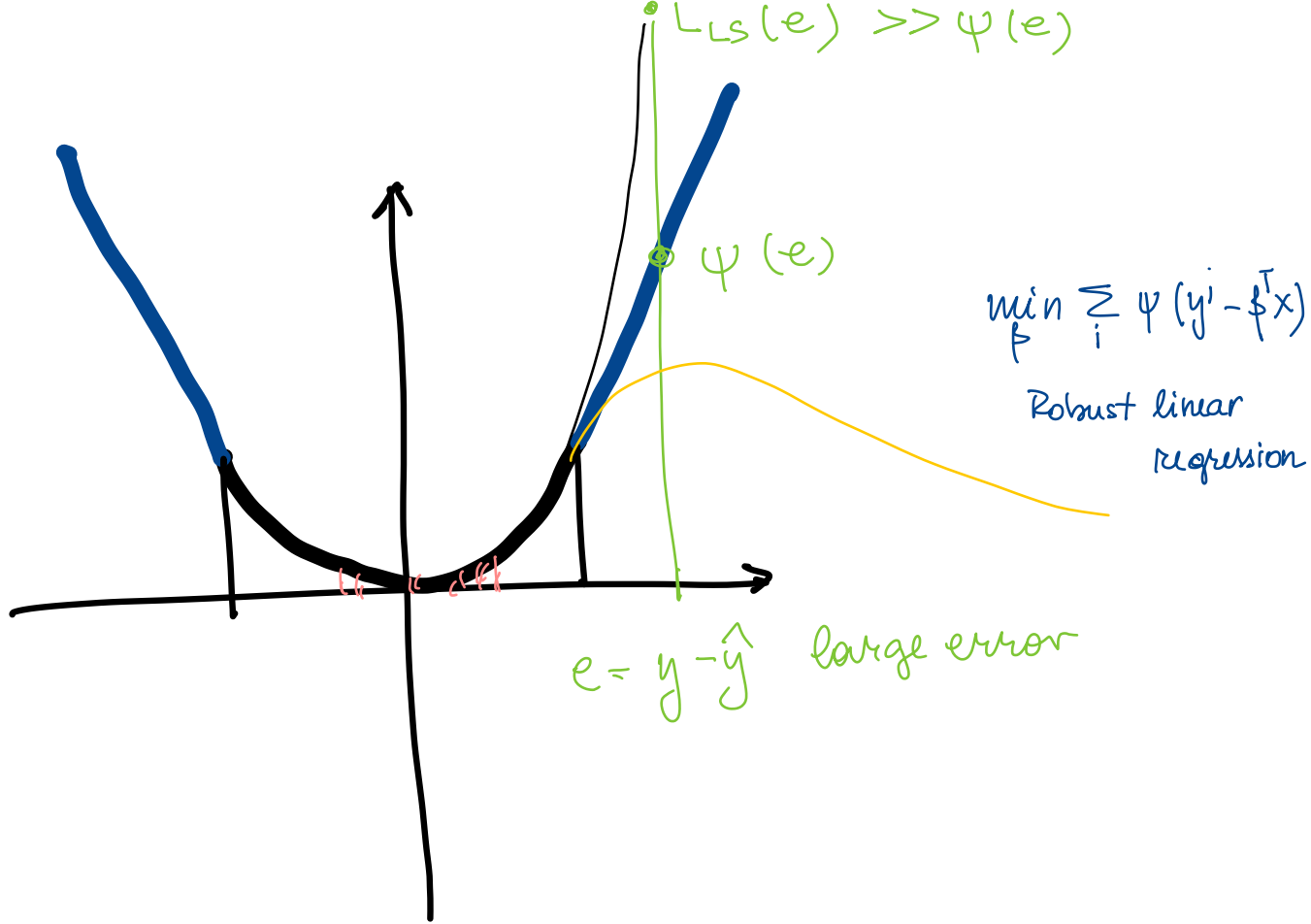  ← "induces sparsity"

  $\|\beta\|_1$

  ↗ Robust

- Huber

  $L_H = \psi(y - \hat{y})$

- −Log Likelihood

  $L_{logL} = -\ln P(y|x)$

  ↳ in model class $\mathcal{F}$

$\psi$

linear

$y - \hat{y}$

$\|y - \hat{y}\|^2$

$L_{LS}(e) \gg \psi(e)$

$\psi(e)$

$\min_{\beta} \sum_i \psi(y^i - \beta^T x)$

Robust linear regression

$e = y - \hat{y}$  large error

# Loss functions for classification

For classification, a natural loss function is the **misclassification error** (also called **0-1 loss**)

$$L_{01}(y, f(x)) = 1_{[y \neq f(x)]} = \begin{cases} 1 & \text{if } y \neq f(x) \quad \text{mistake} \\ 0 & \text{if } y = f(x) \end{cases} \tag{6}$$

Sometimes different errors have different costs. For instance, classifying a HIV+ patient as negative (**a false negative** error) incurs a much higher cost than classifying a normal patient as HIV+ (**false positive** error). This is expressed by **asymmetric misclassification costs**. For instance, assume that a false positive has cost one and a false negative has cost 100. We can express this in the matrix

| $f(x):$ | $+$ | $-$ |
|---|---|---|
| true $:+$ | 0 | 100 ← **miss** |
| $-$ | 1 | 0 |

↑ **false detection / false alarm**

In general, when there are $p$ classes, the matrix $L = [L_{kl}]$ defines the loss, with $L_{kl}$ being the cost of misclassifying as $l$ an example whose true class is $k$.

**Multiclass**

$y \in \{1, \ldots m\}$

$$L = \begin{array}{c|cccc} & 1 & 2 & \cdot & m \\ \hline \text{True } y \; 1 & 0 & & & \\ 2 & & 0 & & \\ \cdot & & & 0 & \\ m & & & & 0 \end{array} \longrightarrow L_{y\hat{y}} = \text{cost of guessing } \hat{y} \text{ when } y \text{ true}$$

# Expected loss and empirical loss

▶ Objective of prediction = to minimize expected loss on future data, i.e.

$$\text{minimize } L(f) = E_{P(X,Y)}[L(Y, f(X))] \text{ over } f \in \mathcal{F} \quad \textbf{model class} \quad (7)$$

*problem*

$f \in \mathcal{F}$

*mature*

We call $L(f)$ above **expected loss**.

## Example (Misclassification error $L_{01}(f)$)

$L_{01}(f) = $ probability of making an error on future data.

$$L_{01}(f) = P[Yf(X) < 0] = E_{P_{XY}}[1_{[Yf(X)<0]}] \quad (8)$$

$$\text{sgn } f(x) \neq y \quad \text{mistake}$$

# Expected loss and empirical loss

▶ Objective of prediction = to minimize expected loss on future data, i.e.

$$\text{minimize } L(f) = E_{P(X,Y)}[L(Y, f(X))] \text{ over } f \in \mathcal{F} \tag{7}$$

We call $L(f)$ above **expected loss**.

▶ $L(f)$ cannot be minimized or even computed directly, because we don't know the data distribution $P_{XY}$.
Therefore, in training predictors, one uses the **empirical** data distribution given by the sample $\mathcal{D}$.

▶ The empirical loss (or **empirical error** or **training error**) is the average loss on $\mathcal{D}$

$$\hat{L}(f) = \frac{1}{n} \sum_{i=1}^{n} L(y^i, f(x^i)) \qquad \hat{L}(f) = \frac{1}{n} \sum_{i=1}^{n} 1_{[y^i f(x^i) < 0]} \tag{8}$$

*averaged over $\mathcal{D}$*

▶ Finally, the value of the **optimal expected loss** for our model class (this is the loss value we are aiming for) is denoted by $L(\mathcal{F})$.

*best $L(f)$ in $\mathcal{F}$*

$$L(\mathcal{F}) = \min_{f \in \mathcal{F}} E_{P(X,Y)}[L(Y, f(X))] \qquad L(f) \tag{9}$$

Note that of all the quantities above, we can only know $\hat{L}(f)$ for a finite number of $f$'s in $\mathcal{F}$.

$\hat{L}(f) =$ empirical Loss

Training ——— $\min_{\mathcal{F}} \hat{L}(f)$,

$\min_{\mathcal{F}} \left[ \hat{L}(f) + \lambda R(f) \right]$

$\llcorner \mathcal{D}$

$\llcorner$ regularization indep of $\mathcal{D}$

other (LDA)

$1 = NN$
linear regression $(L_{LS})$
logistic regression

SVM, Lasso regression
Ridge regression

K - NN
NW

# Bayes loss

problem $\rightarrow$ L
nature $\rightarrow$ $P_{XY}$

▶ How small can the expected loss $L(f)$ be?
It is clear that

$$L(\mathcal{F}) = \min_{f \in \mathcal{F}} L(f) \geq \min_{f} L(f) = L^* \tag{10}$$

where $L^*$ is taken over all possible functions $f$ that take values in $\mathcal{Y}$.

▶ $L^*$ is the absolute minimum loss for the given $P_{XY}$ and it is called the **Bayes loss**.

▶ The Bayes loss is usually not zero

I    $y = f^*(x)$ deterministic $\Rightarrow$ $L^* = 0$

II   $P_{y|x}$ not deterministic $\Rightarrow$ $L^* > 0$

given $x$ :        choose: $\hat{y}$    $E_{P_{y|x}}[L(y, \hat{y})] =$
            $P_{y|x}$ known                    $= \ell(x, \hat{y})$

$\boxed{f^*(x)} = y^*_{(x)} = \underset{\hat{y}}{\text{argmin}}$

# Bayes loss for (binary) classification

▶ Fix $x$ and assume $P_{Y|X}$ known. Then:

- ▶ Label $y$ will have probability $P_{Y|X}(y|x)$ at this $x$.
- ▶ No deterministic guess $f(x)$ for $y$ will make the classification error $E_{P_{Y|X=x}}[L_{01}(y, f(x))]$ (unless $P_{Y|X=x}$ is itself deterministic)
- ▶ Best guess minimizes the probability of being wrong. This is achieved by chosing the most probable class

$$y^*(x) = \underset{y}{\mathrm{argmax}}\, P_{Y|X}(y|x) \qquad \underline{\pm} \qquad (11)$$

*most probable*

- ▶ The probability of being wrong if we choose $y^*(x)$ is $1 - p^*(x)$, where $p^*(x) = \max_y P_{Y|X}(y|x)$.

▶ The **Bayes classifier** is $y^*(x)$ as a function of $x$ and its expected loss is the Bayes loss

$$L_{01}^* = E_{P_X}[1 - p^*(X)] = E_{P_X}[1 - \max_y P[Y|X]] \qquad (12)$$

*Pr[err]*

This shows that the Bayes loss is a property of the problem, via $L$ and $P_{XY}$, and not of any model class or learning algorithm.

## Example

In a classification problem where the class label depends deterministically of the input, the Bayes loss is 0. For example, classifying between written English and written Japanese has (probably) zero Bayes loss.

## Example

Consider the least squares loss and the following data distribution: $P_{Y|X} \sim N(g(X), \sigma^2)$. In other words, the $Y$ values are normally distributed around a deterministic function $g(X)$. In this case, optimal least squares predictor is the mean of $Y$ given $X$, which is equal to $g(X)$. The Bayes loss is the expected squared error around the mean, which is $\sigma^2$. Exercise what is the expression of the Bayes loss if $P_{Y|X} \sim N(g(X), \sigma(X)^2)$?

Exercise What is the Bayes loss if (1) $P(Y|X) \sim N((\beta^*)^T X, \sigma^2 I)$ and the loss is $L_{LS}$; (2) $P(X|Y = \pm 1) \sim N(\mu_{\pm}, \sigma^2 I)$ and the loss is $L_{01}$ (for simplicity, assume $X \in \mathbb{R}$, $\mu_{pm} = \pm 1$, $\sigma = 1$); (3) give a formula for the Bayes loss if we know $P(X|Y = \pm 1), P(Y), Y \in \{\pm 1\}$ and the loss is $L_{01}$. (4) Give an example of a situation when the Bayes loss is 0.

# Bias & Variance for parameter estimation

Pb  $x^{1:n} \sim N(\mu, \sigma^2)$  ← assume true model

Max $\hat{\mu} = ..$

Likelihood $\hat{\sigma}^2 = ...$

$$E[\hat{\mu} - \mu] = 0 \quad \text{r.v. unbiased}$$

$$\mathcal{D}_n \equiv P_X^n \quad n \text{ samples}$$

$$\text{Var} \, \hat{\mu} = \frac{\sigma^2}{n}$$

$$E[\hat{\sigma}^2 - \sigma^2] = -\frac{\sigma^2}{n} \quad \text{biased}$$

$$E[\hat{\sigma}^2] = \sigma^2 \frac{n-1}{n}$$

$$\text{Var} \, \hat{\sigma}^2 = ....$$

Prediction
- estimating $f(x)$  a function !!

  $P_{xy}$  distribution
- no $f^{true}$  (sometimes)