

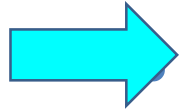
# Lecture 6: Model Selection, Multivariate Classification and Dimensionality Reduction

STAT261: Introduction to Machine Learning

Lecture 6, April 13

# Outline: Multivariate Classification and Dimensionality Reduction

- Multivariate Data and Multiple Measurements
  - Multivariate Parameters
  - Multivariate Gaussian



## Parametric Classification

- Different Covariances
  - Quadratic Discriminant
- Dimensionality Reduction

# Multivariate Data

- Multiple measurements (sensors)
- $d$  inputs/features/attributes:  $d$ -variate
- $N$  instances/observations/examples

$$\mathbf{X} = \begin{bmatrix} X_1^1 & X_2^1 & \dots & X_d^1 \\ X_1^2 & X_2^2 & \dots & X_d^2 \\ \vdots & & & \\ X_1^N & X_2^N & \dots & X_d^N \end{bmatrix}$$

# Multivariate Parameters

$$\text{Mean: } E[\mathbf{x}] = \boldsymbol{\mu} = [\mu_1, \dots, \mu_d]^T$$

$$\text{Covariance: } \sigma_{ij} \equiv \text{Cov}(x_i, x_j)$$

$$\text{Correlation: } \text{Corr}(x_i, x_j) \equiv \rho_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j}$$

$$\Sigma \equiv \text{Cov}(\mathbf{X}) = E[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T] = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1d} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2d} \\ \vdots & & & \\ \sigma_{d1} & \sigma_{d2} & \cdots & \sigma_d^2 \end{bmatrix}$$

# Parameter Estimation

Sample mean  $\mathbf{m}$  :  $m_i = \frac{\sum_{t=1}^N x_i^t}{N}, i = 1, \dots, d$

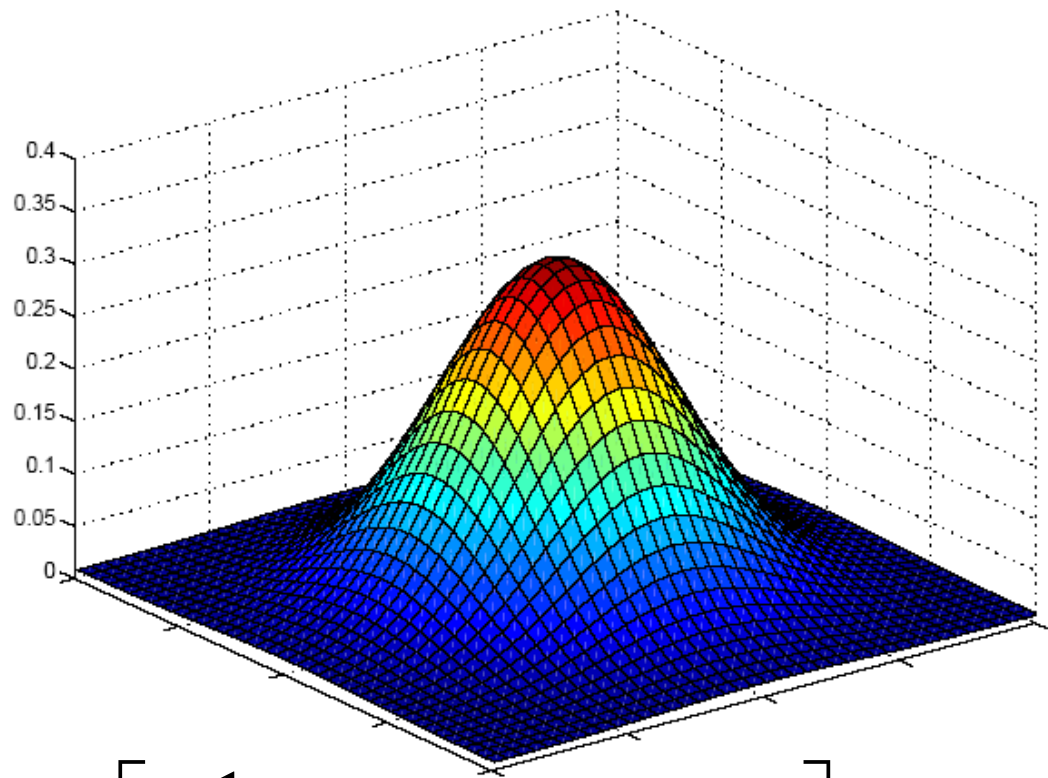
Covariance matrix  $\mathbf{S}$  :  $s_{ij} = \frac{\sum_{t=1}^N (x_i^t - m_i)(x_j^t - m_j)}{N}$

Correlation matrix  $\mathbf{R}$  :  $r_{ij} = \frac{s_{ij}}{s_i s_j}$

# Estimation of Missing Values

- What to do if certain instances have missing attributes?
- Ignore those instances: not a good idea if the sample is small
- Use 'missing' as an attribute: may give information
- Imputation: Fill in the missing value
  - Mean imputation: Use the most likely value (e.g., mean)
  - Imputation by regression: Predict based on other attributes

# Multivariate Normal Distribution



$$\mathbf{x} \sim \mathcal{N}_d(\boldsymbol{\mu}, \Sigma)$$

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

# Multivariate Normal Distribution

- Mahalanobis distance:  $(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$   
measures the distance from  $\mathbf{x}$  to  $\boldsymbol{\mu}$  in terms of  $\boldsymbol{\Sigma}$   
(normalizes for difference in variances and correlations)
- Bivariate:  $d = 2$

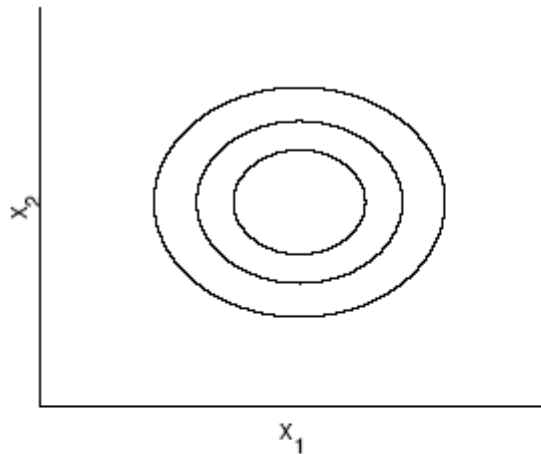
$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$$

$$p(x_1, x_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left[-\frac{1}{2(1-\rho^2)}(z_1^2 - 2\rho z_1 z_2 + z_2^2)\right]$$
$$z_i = (x_i - \mu_i) / \sigma_i$$

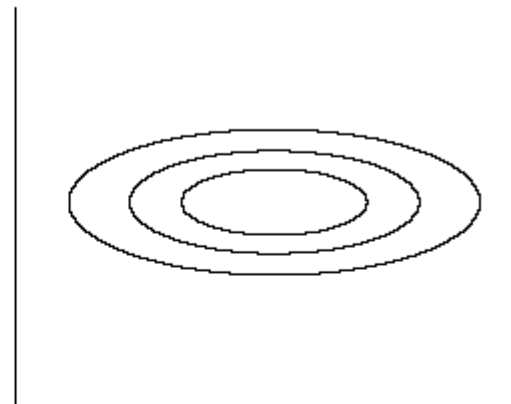


# Bivariate Normal

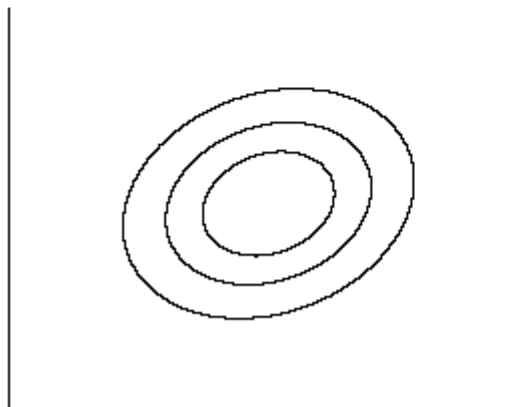
$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) = \text{Var}(x_2)$



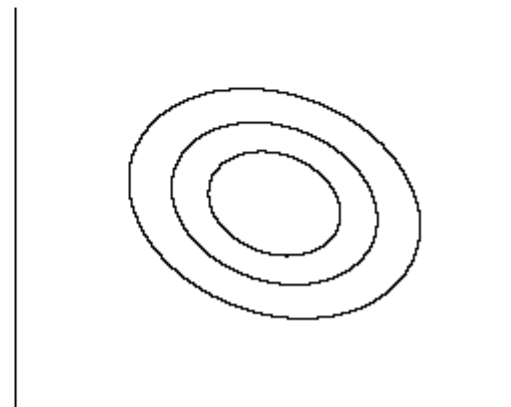
$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) > \text{Var}(x_2)$



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

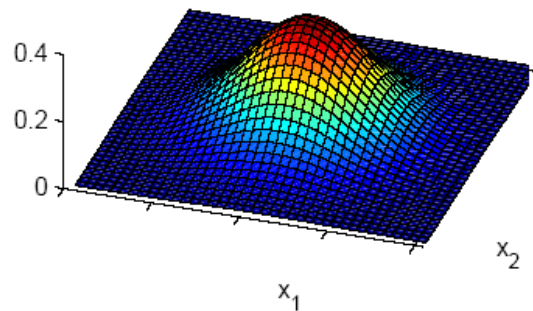


$\text{Cov}(x_1, x_2) < 0$

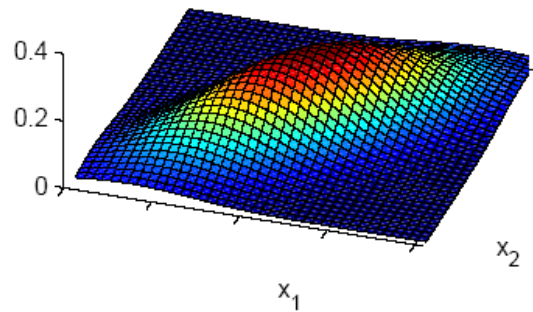


# Bivariate Normal

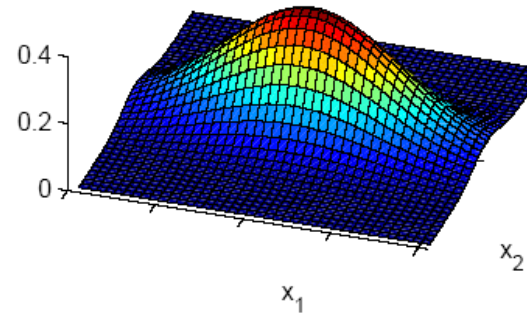
$$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) = \text{Var}(x_2)$$



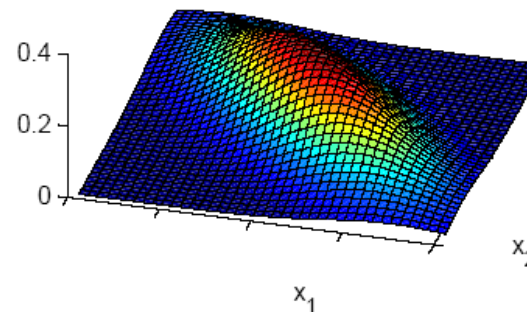
$$\text{Cov}(x_1, x_2) > 0$$



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



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



# Independent Inputs: Naive Bayes

- If  $x_i$  are independent, off diagonals of  $\Sigma$  are 0, Mahalanobis distance reduces to weighted (by  $1/\sigma_i$ ) Euclidean distance:

$$p(\mathbf{x}) = \prod_{i=1}^d p_i(x_i) = \frac{1}{(2\pi)^{d/2} \prod_{i=1}^d \sigma_i} \exp \left[ -\frac{1}{2} \sum_{i=1}^d \left( \frac{x_i - \mu_i}{\sigma_i} \right)^2 \right]$$

- If variances are also equal, reduces to Euclidean distance

# Parametric Classification

Classification: Maximum likelihood

Pick class that would make observation most likely

- If  $p(\mathbf{x} | C_i) \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$

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

- Discriminant functions

$$\begin{aligned} g_i(\mathbf{x}) &= \log p(\mathbf{x} | C_i) + \log P(C_i) \\ &= -\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\boldsymbol{\Sigma}_i| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \log P(C_i) \end{aligned}$$

# Estimation of Parameters

$$\hat{P}(C_i) = \frac{\sum_t r_i^t}{N}$$

$$\mathbf{m}_i = \frac{\sum_t r_i^t \mathbf{x}^t}{\sum_t r_i^t}$$

$$\mathbf{S}_i = \frac{\sum_t r_i^t (\mathbf{x}^t - \mathbf{m}_i)(\mathbf{x}^t - \mathbf{m}_i)^T}{\sum_t r_i^t}$$

$$g_i(\mathbf{x}) = -\frac{1}{2} \log |\mathbf{S}_i| - \frac{1}{2} (\mathbf{x} - \mathbf{m}_i)^T \mathbf{S}_i^{-1} (\mathbf{x} - \mathbf{m}_i) + \log \hat{P}(C_i)$$

# Parametric Classification: Different Variances

- Different  $\mathbf{S}_i$
- Quadratic discriminant

$$g_i(\mathbf{x}) = -\frac{1}{2} \log |\mathbf{S}_i| - \frac{1}{2} (\mathbf{x}^T \mathbf{S}_i^{-1} \mathbf{x} - 2 \mathbf{x}^T \mathbf{S}_i^{-1} \mathbf{m}_i + \mathbf{m}_i^T \mathbf{S}_i^{-1} \mathbf{m}_i) + \log \hat{P}(C_i)$$

$$= \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

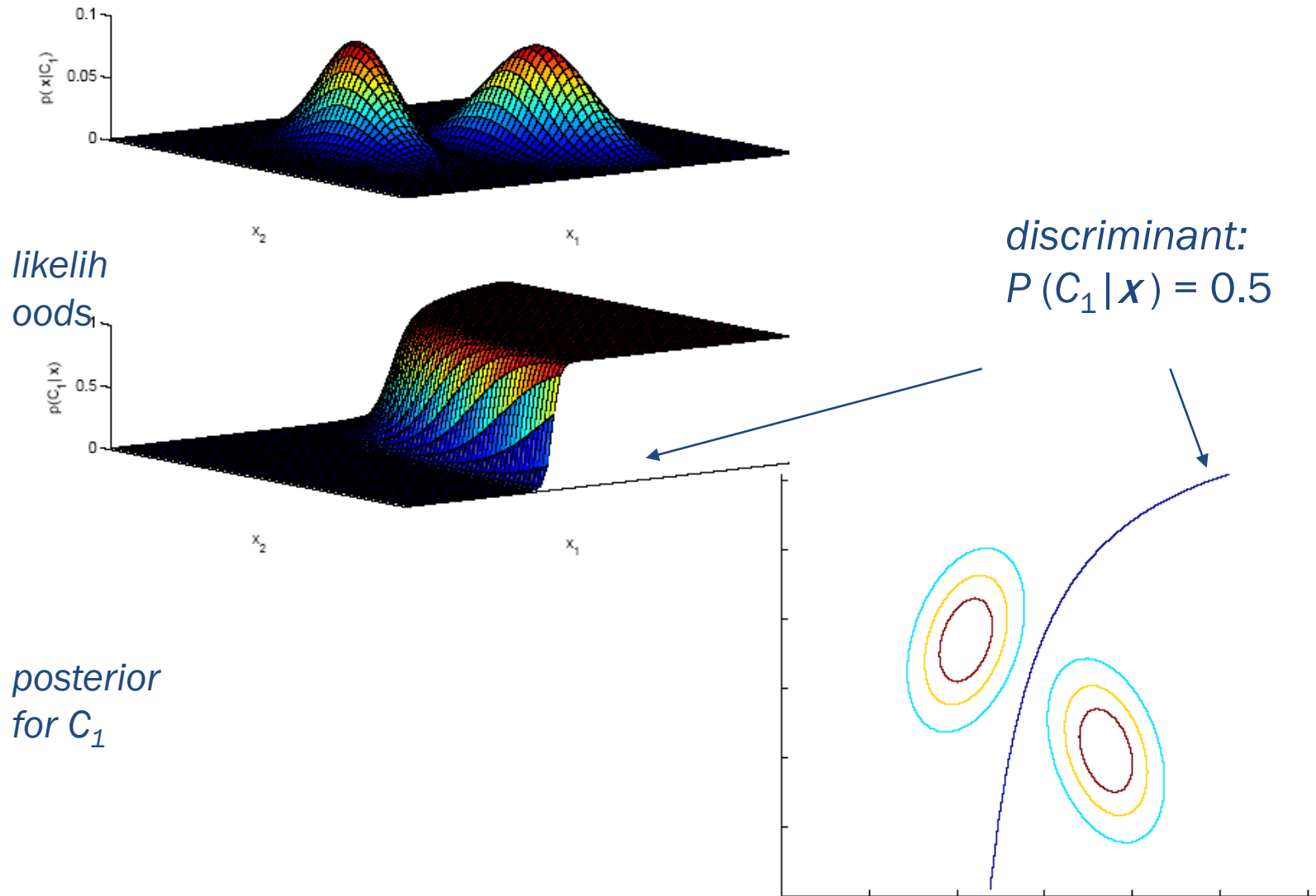
where

$$\mathbf{W}_i = -\frac{1}{2} \mathbf{S}_i^{-1}$$

$$\mathbf{w}_i = \mathbf{S}_i^{-1} \mathbf{m}_i$$

$$w_{i0} = -\frac{1}{2} \mathbf{m}_i^T \mathbf{S}_i^{-1} \mathbf{m}_i - \frac{1}{2} \log |\mathbf{S}_i| + \log \hat{P}(C_i)$$

# Parametric Classification: Different Variances



# Common Covariance Matrix $\mathbf{S}$

- Shared common sample covariance  $\mathbf{S}$

- Discriminant reduces to 
$$\mathbf{S} = \sum_i \hat{P}(C_i) \mathbf{S}_i$$

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mathbf{m}_i)^T \mathbf{S}^{-1}(\mathbf{x} - \mathbf{m}_i) + \log \hat{P}(C_i)$$

which is a linear discriminant

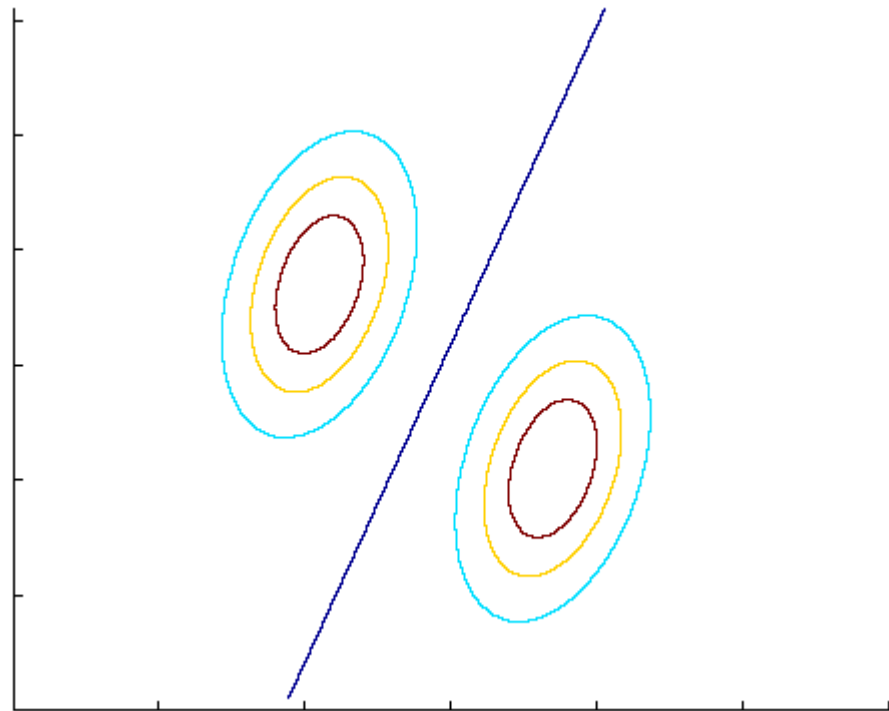
$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

where

$$\mathbf{w}_i = \mathbf{S}^{-1} \mathbf{m}_i \quad w_{i0} = -\frac{1}{2} \mathbf{m}_i^T \mathbf{S}^{-1} \mathbf{m}_i + \log \hat{P}(C_i)$$



# Common Covariance Matrix $\mathbf{S}$



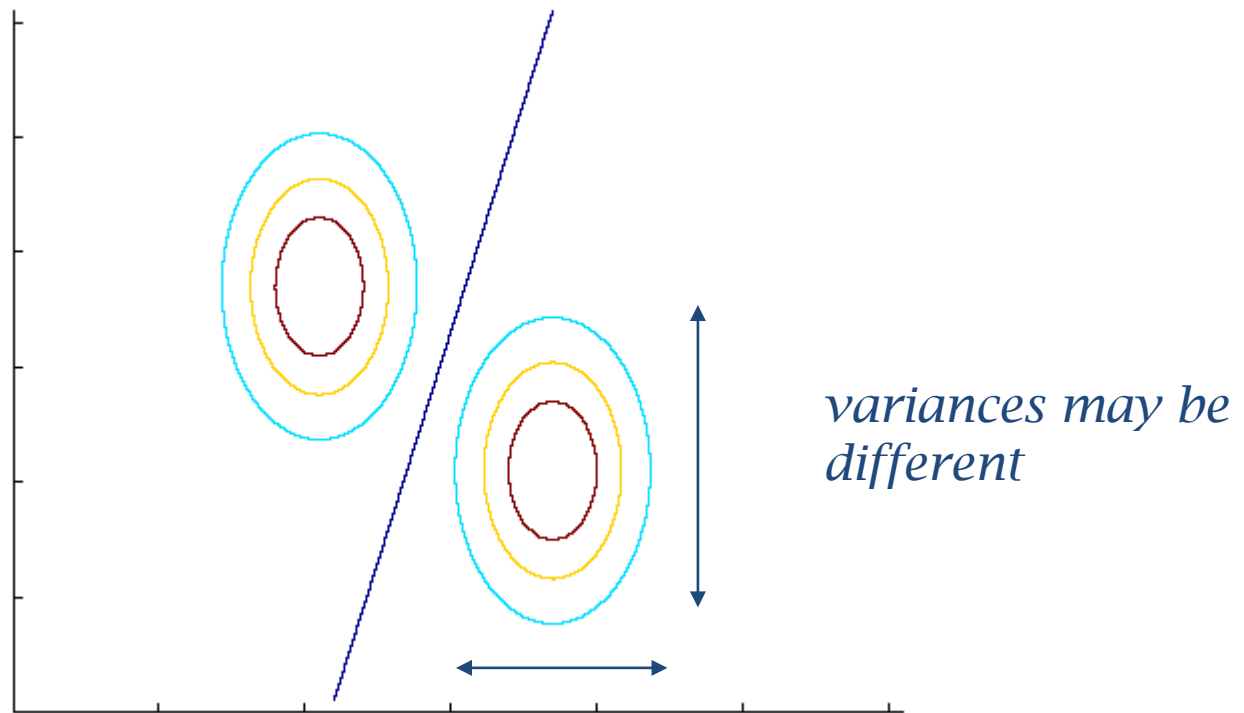
# Classification: Diagonal $\mathbf{S}$ 's

- When  $x_j, j = 1, \dots, d$ , are independent,
- $\Sigma$  is diagonal
- $p(\mathbf{x} | C_i) = \prod_j p(x_j | C_i)$  (Naive Bayes' assumption)

$$g_i(\mathbf{x}) = -\frac{1}{2} \sum_{j=1}^d \left( \frac{x_j^t - m_{ij}}{s_j} \right)^2 + \log \hat{P}(C_i)$$

- Classify based on weighted Euclidean distance  
(in  $s_j$  units) to the nearest mean

# Classification: Diagonal $\mathbf{S}$ matrices



# Classification: Diagonal $\mathbf{S}$ , equal variances

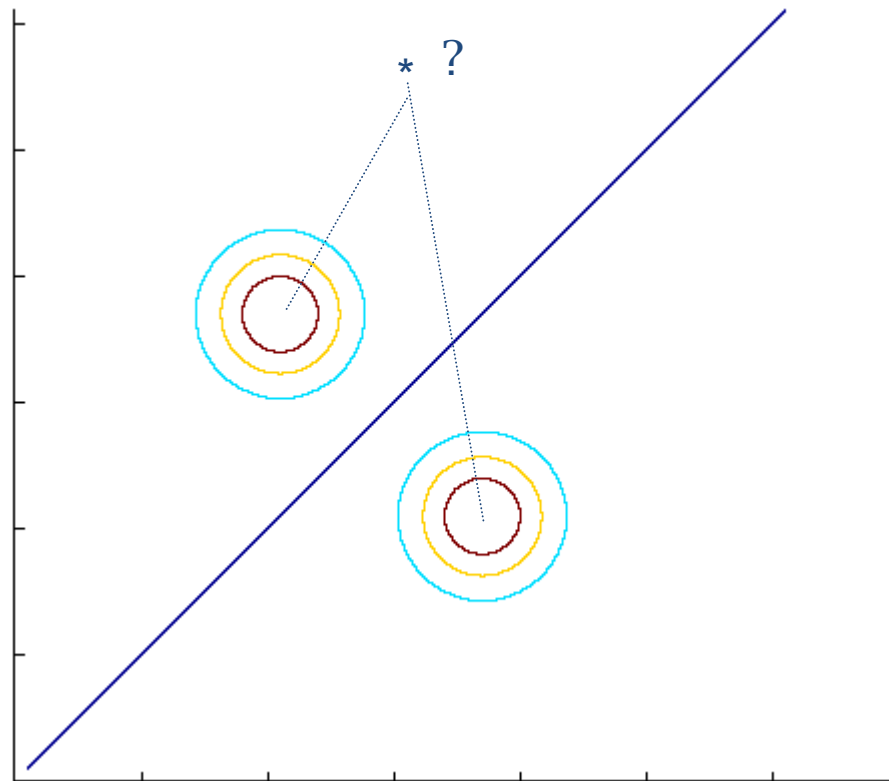
- Nearest mean classifier:

Classify based on Euclidean distance to the nearest mean

$$\begin{aligned} g_i(\mathbf{x}) &= -\frac{\|\mathbf{x} - \mathbf{m}_i\|^2}{2s^2} + \log \hat{P}(C_i) \\ &= -\frac{1}{2s^2} \sum_{j=1}^d (x_j^t - m_{ij})^2 + \log \hat{P}(C_i) \end{aligned}$$

- Each mean can be considered a prototype or template  
This is sometimes called template matching

# Diagonal $\mathbf{S}$ , equal variances



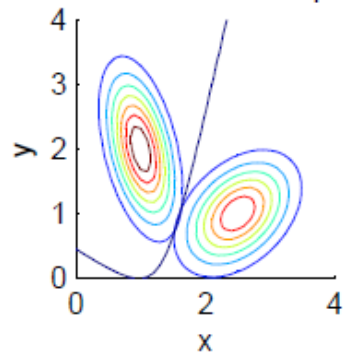
# Model Selection

| <i>Assumption</i>           | <i>Covariance matrix</i>                        | <i>No of parameters</i> |
|-----------------------------|---|-------------------------|
| Shared, Hyperspheric        | $\mathbf{S}_i = \mathbf{S} = s^2 \mathbf{I}$    | 1                       |
| Shared, Axis-aligned        | $\mathbf{S}_i = \mathbf{S}$ , with $s_{ij} = 0$ | $d$                     |
| Shared, Hyperellipsoidal    | $\mathbf{S}_i = \mathbf{S}$                     | $d(d+1)/2$              |
| Different, Hyperellipsoidal | $\mathbf{S}_i$                                  | $K d(d+1)/2$            |

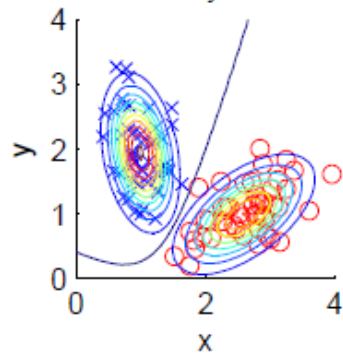
- As we increase complexity (less restricted  $\mathbf{S}$ ), bias decreases and variance increases
- Assume simple models (allow some bias) to control variance (regularization)

# Population Likelihoods and Posteriors

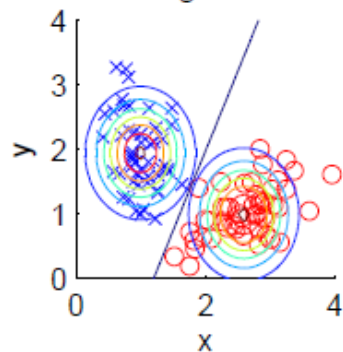
Population likelihoods and posteriors



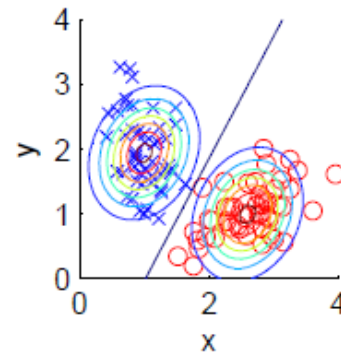
Arbitrary covar.



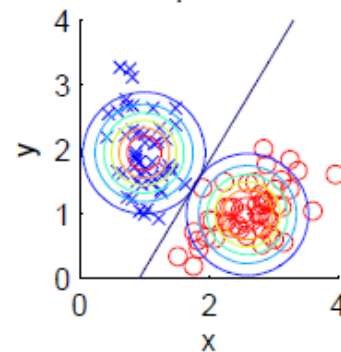
Diag. covar.



Shared covar.



Equal var.



# Classification: Discrete Features

- Binary features:  $p_{ij} \equiv p(x_j=1 | C_i)$   
if  $x_j$  are independent (Naive Bayes')

$$p(\mathbf{x} | C_i) = \prod_{j=1}^d p_{ij}^{x_j} (1 - p_{ij})^{(1-x_j)}$$

the discriminant is linear

$$\begin{aligned} g_i(\mathbf{x}) &= \log p(\mathbf{x} | C_i) + \log P(C_i) \\ &= \sum_j [x_j \log p_{ij} + (1 - x_j) \log (1 - p_{ij})] + \log P(C_i) \end{aligned}$$

Estimated parameters  $\hat{p}_{ij} = \frac{\sum_t x_j^t r_i^t}{\sum_t r_i^t}$



# Discrete Features

- Multinomial (1-of- $n_j$ ) features:  $x_j \in \{v_1, v_2, \dots, v_{n_j}\}$

$$p_{ijk} \equiv p(z_{jk}=1 | C_i) = p(x_j = v_k | C_i)$$

if  $x_j$  are independent

$$p(\mathbf{x} | C_i) = \prod_{j=1}^d \prod_{k=1}^{n_j} p_{ijk}^{z_{jk}}$$

$$g_i(\mathbf{x}) = \sum_j \sum_k z_{jk} \log p_{ijk} + \log P(C_i)$$

$$\hat{p}_{ijk} = \frac{\sum_t z_{jk}^t r_i^t}{\sum_t r_i^t}$$

(6)

## The Curse of Dimensionality

The examples of Bayes Decision theory are misleading because they are given in low-dimensional spaces (1-dim, or 2-dim)

Many pattern classification tasks occur in high dimensional spaces. In these spaces our geometric intuitions are often wrong.

EG. Consider the volume of a sphere of radius  $r=1$  in  $D$  dimensions.

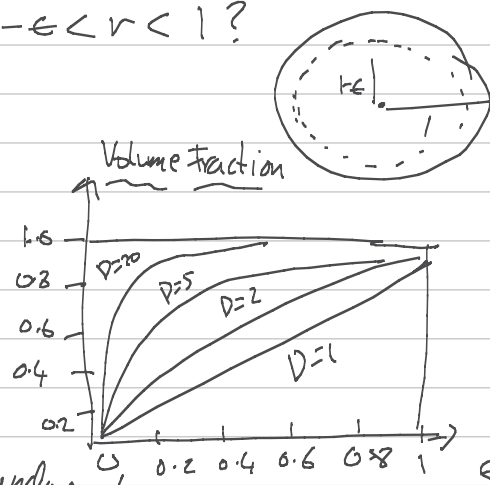
What fraction of its volume lies in the region between  $1-\epsilon < r < 1$ ?

$$V_D(r) = K_D r^D$$

$$\frac{V_D(1) - V_D(1-\epsilon)}{V_D(1)} = 1 - (1-\epsilon)^D$$

For large  $D$ , the volume fraction tends to 1 even for small  $\epsilon$ .

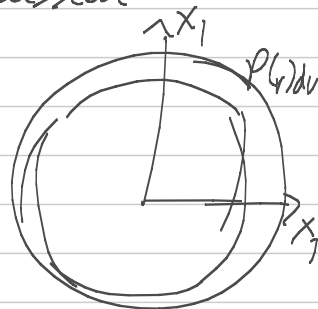
Most of the volume is at the boundary!



(7)

e.g. Behaviour of a Gaussian distribution.

In 1-D,  $p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/2\sigma^2}$



In 2-D

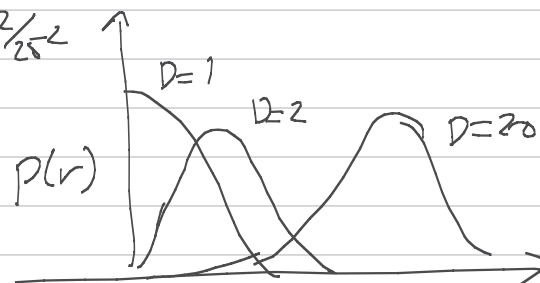
$$p(x_1, x_2) = \frac{1}{2\pi\sigma^2} e^{-\frac{(x_1^2 + x_2^2)}{2\sigma^2}}$$

Let  $r = \sqrt{x_1^2 + x_2^2}$ .

Then  $p(r) = \frac{r}{2\pi\sigma^2} e^{-r^2/2\sigma^2}$

In higher dimension

$$p(r) = \frac{r^{D-1}}{\Gamma} e^{-r^2/2\sigma^2}$$



So in high dimensions

most of the probability

mass of the Gaussian is

concentrated on a thin shell away from the center of the Gaussian.

(8)

Learning probability distributions  
in high dimensions can require a lot of data.

E.G. Gaussian Distribution in  $D$  dimensions.

mean -  $\underline{\mu}$   $D$  dimensions.

Covariance -  $\underline{\Sigma}$   $\frac{D(D+1)}{2}$  Dimensions.

This is  $O(D^2)$ , not too bad.

But suppose we represent the data by  
a histogram with  $B$  bins per dimension.

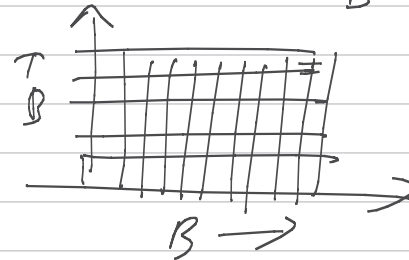
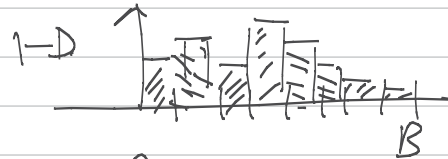
$B$  bins in  $D=1$

$B^2$  bins in  $D=2$

$B^D$  bins in  $D$  dimensions.

Exponential growth!

Requires exponential amount  
of data to learn the distribution.

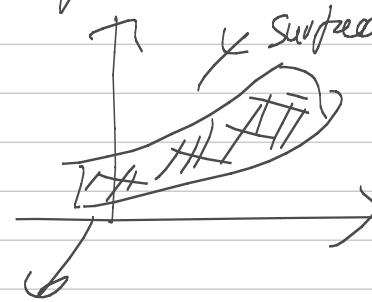


(9)

How to deal with the curse of dimensionality?

In practice, data typically lies on some low-dimensional surface in the high dimensional space.

So the effective dimension of the data may be a lot smaller than the dimension of the space.



(\*) Dimension Reduction Methods attempt to reduce the dimension by seeking this low dimensional surface. (Not always easy).

(\*) Modeling, if we can guess distributions for the data (e.g. Gaussian) then the dependence on the dimension is not too bad.

(\*) Concentrate on the Decision Boundary ~ there may be enough data to learn the decision boundary even if we cannot learn the distributions.

## (1b) Bias and Variance

This is a classical statistics perspective on generalization. First we need to introduce some statistics terminology.

Suppose we want to estimate a continuous quantity  $\theta$  - e.g. the mean/variance of a Gaussian distribution (more about this in the next lecture), or the parameters of a regression line (see below) - then statisticians use an estimator.

The estimator is based on a set  $X_N = \langle x_i : i=1 \dots N \rangle$  of examples - drawn from an unknown distribution  $P(x)$ .  
i.i.d.  $P(X_N) = \prod_{i=1}^N P(x_i)$

The task is to estimate a property  $\theta$  by an estimator  $\hat{\theta} = g(X_N)$ . E.g. like a classification rule - but based on the set  $X_N$  and  $\theta$  is continuous.

For example: let  $\theta = (\mu, \sigma)$  be the mean and variance of the data (data is one-dimensional in this example).

$$\text{then } X_N = \langle x_i : i=1 \dots N \rangle \\ \hat{\mu}(X_N) = \frac{1}{N} \sum_{i=1}^N x_i, \quad \hat{\sigma}^2(X_N) = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu}(X_N))^2$$

$$g(X_N) = (\hat{\mu}(X_N), \hat{\sigma}^2(X_N))$$

Note: that the estimator is a function of the set  $X_N$ , this will be important later.

(11)

To evaluate the estimator,  
we want to measure how much it differs from  
the correct  $\theta$ . It is attractive to use a  
quadratic errors (this helps the analysis)

$$(g(X_n) - \theta)^2 \quad \text{— but this depends on the data set } X_n$$

So we need to get the expected error with  
respect to the set  $X_n$ ,

$$r(g, \theta) = E_X [g(X) - \theta]^2 = \int (g(x) - \theta)^2 P(x) dx$$

mean square error. distribution over set  $X = \{x_i\}$ .

$$b_\theta(g) = E_X [g(X)] - \theta, \quad \text{bias of estimator}$$

If  $b_\theta(g) = 0$  for all  $\theta$ , then  $g(\cdot)$  is an unbiased estimator of  $\theta$

Eg. consider  $\hat{\mu}(X) = \frac{1}{N} \sum_{i=1}^N x_i$ , estimate of the mean

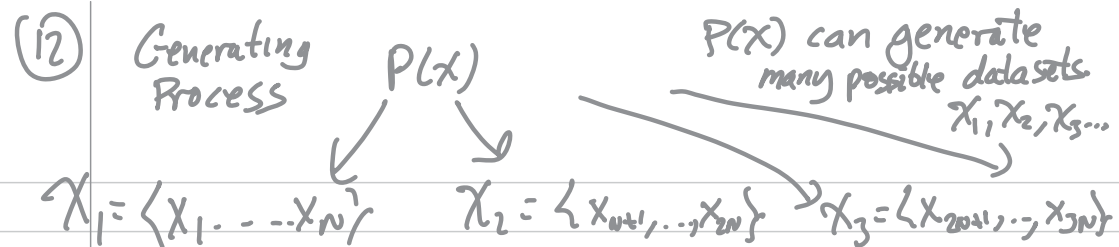
$$E_X [\hat{\mu}(X)] = E_X \left[ \frac{1}{N} \sum_{i=1}^N x_i \right] = \frac{1}{N} \sum_{i=1}^N E_{x_i}(x_i)$$
$$= \frac{1}{N} N \mu = \mu$$

$\sum_{i=1}^N x_i P(x_i) = \mu$   
mean of the distribution that generated the data.

hence  $\hat{\mu}(\cdot)$  is an unbiased estimator of  $\mu$ .

( $\hat{\mu}(X)$  will depend on  $X$ , but  
on average it gives you the right answer).

We can compute the variance of the  
estimator — i.e. how much it varies depending on  $X$ .



$g(X_1), g(X_2), g(X_3)$ , these estimators will vary  
 Calculate their mean  $\rightarrow$  e.g.  $E_X\{g(X)\}$

Calculate the variance  $\rightarrow \text{Var}(g(X)) = E_X\{(g(X) - E_X\{g(X)\})^2\}$

For the estimator  $\hat{\mu}(X)$  of the mean,  
 we know that  $E_X\{\hat{\mu}(X)\} = \mu$  unbiased (previous page)

$$\begin{aligned} \text{Var}_X(\hat{\mu}) &= \text{Var}_X\left(\sum_{i=1}^N X_i\right) = \frac{1}{N^2} \sum_{i=1}^N \text{Var}_{X_i}(X_i) \\ &= \frac{N}{N^2} \sigma^2 = \frac{\sigma^2}{N} \end{aligned}$$

$\underbrace{\sum_{i=1}^N \text{Var}_{X_i}(X_i)}_{\sum_{i=1}^N P(X_i)(X_i - \mu)^2}$

Hence, the variance of the estimator tends to zero as  $N \rightarrow \infty$  with fall off rate  $1/N$

Note: this  $1/N$  fall-off rate is true for any linear estimator - i.e.  $g(X)$  is a linear function of the elements  $X_1, \dots, X_N$  in set  $X$ .

Next, consider the bias of  $\hat{\sigma}^2(X) = \frac{1}{N} \sum_{i=1}^N (X_i - \hat{\mu}_i)^2$

By similar analysis (to analysis of  $\hat{\mu}(X)$ ) we find that  
 $E_X[\hat{\sigma}^2(X)] = \frac{1}{N} \left\{ \sum_{i=1}^N E_{X_i}(X_i^2) - N \mu^2 \right\} = \frac{N-1}{N} E_X(X_i^2) = \frac{N-1}{N} (\mu^2 + \sigma^2)$   
 (uses  $\sum_{i=1}^N (X_i - \mu)^2 = \sum_{i=1}^N X_i^2 - N\mu^2$ ). biased, but asymptotically unbiased as  $N \rightarrow \infty$



### (13) Bias Variance Dilemma:

Note Title

3/29/2008

Dataset  $\mathcal{X} = \langle (x_i, y_i) : i = 1 \text{ to } N \rangle$ ,  
sampled from  $P(x, y) = p(y|x)P(x)$   
 $y$  is continuous valued.

Let  $g(x)$  be an estimator of  $y$

Claim 1:  $\underbrace{\langle (y - g(x))^2 \rangle_{P(y|x)}}_{\text{Expected error w.r.t. } P(y|x)} = \underbrace{(y - \langle y \rangle_{P(y|x)})^2}_{\text{Variance of the process}} + \underbrace{(\langle y \rangle_{P(y|x)} - g(x))^2}_{\text{Squared error}}$

Here  $\langle f(y, x) \rangle_{P(y|x)}$   
 $= \sum_y f(y, x) P(y|x)$

→ Nothing we can do  
about this - indep of  
our estimator  $g(\cdot)$

Proof: Write  $(y - g(x))^2 = (y - \langle y \rangle_{P(y|x)} + \langle y \rangle_{P(y|x)} - g(x))^2$   
 $= (y - \langle y \rangle_{P(y|x)})^2 + (\langle y \rangle_{P(y|x)} - g(x))^2 + 2(y - \langle y \rangle_{P(y|x)})(\langle y \rangle_{P(y|x)} - g(x))$

Take expectation w.r.t.  $P(y|x)$

$$\text{Then } \langle (y - g(x))^2 \rangle_{P(y|x)} = \langle (y - \langle y \rangle_{P(y|x)})^2 \rangle + (\langle y \rangle_{P(y|x)} - g(x))^2$$

because  $(\langle y \rangle_{P(y|x)} - g(x))^2$  is indep of  $y$  and  $\langle (y - \langle y \rangle_{P(y|x)})^2 \rangle_{P(y|x)} = 0$

Claim 1 says we can decompose the expected error  
(w.r.t.  $P(y|x)$ ) into part we have no control over  $\langle (y - \langle y \rangle_{P(y|x)})^2 \rangle$   
and part which depends on  $g(\cdot)$  and the data  $\mathcal{X}$ .  
 $(\langle y \rangle_{P(y|x)} - g(x))^2$

(14) Next we study the expectation

of  $(\langle y \rangle_{p(y|x)} - g(x))^2$  with respect to  $P(X)$

i.e. how it depends on the particular sample  $X = (x_1, \dots, x_n)$  from  $P(X)$ .

Claim II  $\langle (\langle y \rangle_{p(y|x)} - g(x))^2 \rangle_{p(x)} = (\langle y \rangle_{p(y|x)} - \langle g(x) \rangle_{p(x)})^2 + \langle (g(x) - \langle g(x) \rangle_{p(x)})^2 \rangle_{p(x)}$

The Bias-Variance Result.

Bias  $\swarrow$   
Variance  $\searrow$

The first term depends on the bias - the difference between the best estimate  $\langle y \rangle_{p(y|x)}$  (if we knew the distribution) and the expectation of our estimator  $\langle g(x) \rangle_{p(x)}$

The second term is the variance of the estimator  $g(x)$  - i.e. how much it depends on the sample set  $x$  (noisy)

Proof: Write  $(\langle y \rangle_{p(y|x)} - g(x))^2$

$$= (\langle y \rangle_{p(y|x)} - \langle g(x) \rangle_{p(x)} + \langle g(x) \rangle_{p(x)} - g(x))^2$$
$$= (\langle y \rangle_{p(y|x)} - \langle g(x) \rangle_{p(x)})^2 + (\langle g(x) \rangle_{p(x)} - g(x))^2 + 2(\langle y \rangle_{p(y|x)} - \langle g(x) \rangle_{p(x)})(\langle g(x) \rangle_{p(x)} - g(x))$$

Take expectation wrt.  $P(x)$

$$\langle (\langle y \rangle_{p(y|x)} - g(x))^2 \rangle_{p(x)} = (\langle y \rangle_{p(y|x)} - \langle g(x) \rangle_{p(x)})^2 + \langle (g(x) - \langle g(x) \rangle_{p(x)})^2 \rangle_{p(x)}$$

Because  $(\langle y \rangle_{p(y|x)} - g(x))_{p(x)}$  is independent of  $x$  and  $\langle g(x) - \langle g(x) \rangle_{p(x)} \rangle_{p(x)} = 0$

(15)

What does this mean?

Distribution  $P(x)$

Dataset  $X_1 = \{x_1, \dots, x_N\}$ ,  $X_2 = \{x_{N+1}, \dots, x_{2N}\}$   
 $\dots$   $X_m = \{x_{(m-1)N+1}, \dots, x_{mN}\}$

For each data we get an estimate of  $y$

$g(x_1), g(x_2) \dots g(x_m)$

The mean estimate is  $\bar{g} = \frac{1}{m} \sum_{i=1}^m g(x_i)$

the variance is  $\text{Var}(\bar{g}) = \frac{1}{m} \sum_{i=1}^m (g(x_i) - \bar{g})^2$

To get good generalization we want the variance to be small, so that it isn't sensitive to the data we have trained the classifier on.

Ideally we want to have a classifier  $g(\cdot)$  which has small bias and variance.

In practice, there is often a trade-off between bias and variance.

A complex classifier can give a good fit to the data (compared to a simple classifier) but can have high variance because it over-fits the data. So it gives different results on different datasets.



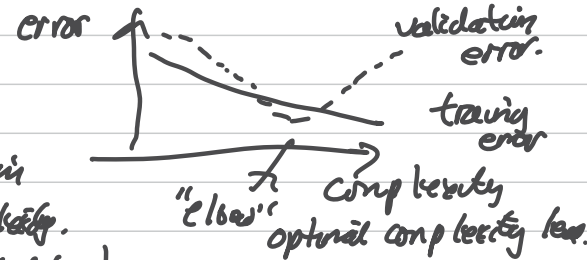
(17)

What to do? How to test generalization?

Cross-validation → divide dataset into two parts as training & validation set.

Train models of different complexity and test their error on the validation set.

As model complexity increases, error on training set decreases. But error on validation set decreases then increases.



regularization

augmented error function

$\tilde{E} = \text{error on data} + \lambda \cdot \text{model complexity}$

( $\lambda$  optimized using cross-validation)

structural risk minimization (Vapnik) (see dim)

— also penalizes model complexity.

Minimum Description Length (Rissanen)

penalize complexity by cost of encoding model.

Bayesian Model Selection. if some prior knowledge

$$P(\text{model} | \text{data}) = \frac{P(\text{data} | \text{model}) P(\text{model})}{P(\text{data})}$$

(gives higher prob. to simpler models)