# Maximum Likelihood
# and
# Expectation
# Maximization

# Probability Basics

Ex

Ex

Sample Space

$\{1, \cdots, 6\}$

$\mathbb{R}$

Distribution

$P_1 \cdots P_6 \geq 0, \Sigma P_i = 1$

p.d.f $f(x) \geq 0, \int_{\mathbb{R}} f(x) dx = 1$

eg $P_1 = \cdots = P_6 = \frac{1}{6}$

$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

Population vs Sample

Population mean

$\mu = \Sigma i P_i$

$\mu = \int x f(x) dx$

Population Variance

$\sigma^2 = \Sigma (i - \mu)^2 P_i$

$\sigma^2 = \int (x-\mu)^2 f(x) dx$

Sample mean

$\Sigma x_i / n$

Sample Variance

$\Sigma (x_i - \bar{x})^2 / n$

1

# Parameter Estimation

Assuming sample $x_1 \ldots x_n$ is from parametric distribution $f(x \mid \Theta)$, estimate $\Theta$.

E.g. $f(x) = \dfrac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

$$\Theta = (\mu, \sigma^2)$$

# Maximimum Likelihood Estimation

One (of many) approaches to parameter est.

Likelihood of $x_1 \cdots x_n =$

$$L(x_1, \ldots, x_n | \theta) = \prod_{i=1}^{n} f(x_i | \theta)$$

assume indp.

View this as a function of $\theta$
what $\theta$ maximizes the likelihood

Typical approach: $\frac{\partial}{\partial \theta} L(x | \theta) = 0$

or $\frac{\partial}{\partial \theta} \ln L(x | \theta) = 0$

# Example

$X_1 \cdots X_n$ coin flips; $\theta =$ prob of heads

$n_0$ tails, $n_1$ heads, $n_0 + n_1 = n$

$$L(X_1 \cdots X_n | \theta) = (1-\theta)^{n_0} (\theta)^{n_1}$$

$$\ln L = n_0 \ln(1-\theta) + n_1 \ln\theta$$

$$\frac{d}{d\theta} \ln L = \frac{-n_0}{1-\theta} + \frac{n_1}{\theta} = 0$$

$$n_0 \theta = n_1(1-\theta)$$

$$(n_0 + n_1)\theta = n_1$$

$$\boxed{\theta = \frac{n_1}{n}}$$

# Example

$$x_i \sim N(\mu, \sigma^2), \quad \sigma^2 = 1, \quad \mu \text{ unknown}$$

$$L(x_1 .. x_n | \theta) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}} e^{-\frac{(x_i - \theta)^2}{2}}$$

$$\ln L(x_1 .. x_n | \theta) = \sum -\frac{1}{2} \ln 2\pi - \frac{(x_i - \theta)^2}{2}$$

$$\frac{dL}{d\theta} = -\sum (x_i - \theta) = 0$$

$$\left( \sum x_i \right) - n\theta = 0$$

$$\boxed{\theta = \sum x_i / n}$$

Example $\quad X_i \sim N(\mu, \sigma^2)$ , both unknown

$$\ln L(x_1 \cdots x_n | \Theta_1, \Theta_2) = \sum -\frac{1}{2} \ln 2\pi \Theta_2 - \frac{(x_i - \Theta_1)^2}{2\Theta_2}$$

$$\frac{\partial}{\partial \Theta_1} \ln L(x_1 \cdots x_n | \Theta_1 \Theta_2) = \sum \frac{x_i - \Theta_1}{\Theta_2} = 0$$

$$\Rightarrow \boxed{\Theta_1 = \sum x_i / n}$$

$$\frac{\partial}{\partial \Theta_2} \ln L(x_1 \cdots x_n | \Theta_1 \Theta_2) = \sum -\frac{2\pi}{2 \cdot 2\pi \Theta_2} + \frac{(x_i - \Theta_1)^2}{2\Theta_2^2}$$

$$\sum -\frac{1}{2\theta_2} + \frac{(x_i - \theta_1)^2}{2\theta_2^2} = 0$$

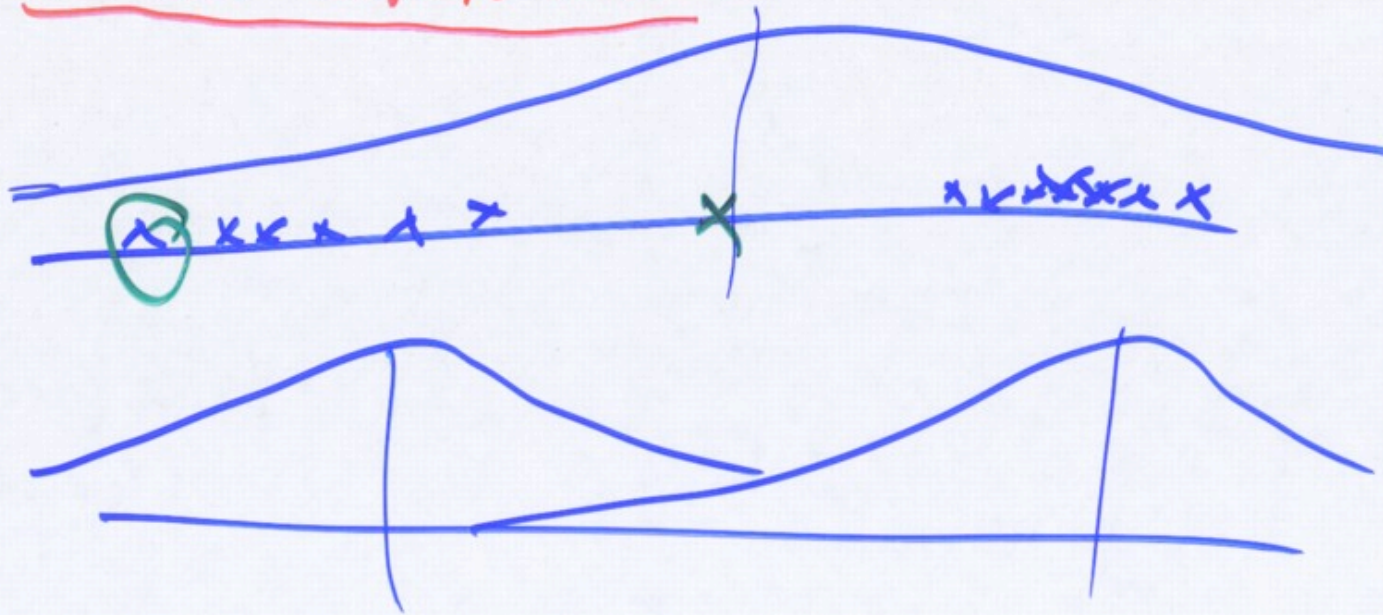$$\sum (x_i - \theta_1)^2 = n\theta_2$$

$$\boxed{\theta_2 = \sum (x_i - \theta_1)^2 / n}$$

A Biased (but consistent) estimate of population variance

An Example of __Overfitting__

Unbiased estimate: $\displaystyle\sum_{i=1}^{n} \frac{(x_i - \theta_1)^2}{n - 1}$

# A More Complex Problem



2 distributions $\quad f_1(x)$ , $f_2(x)$

$\qquad\qquad\qquad f_1(x|\theta_1)$ , $f_2(x|\theta_2)$

Mixing parameter $\quad \tau_1 \qquad \tau_2$

$\qquad\qquad\qquad\qquad \tau_1 + \tau_2 = 1$

Likelihood

$$L(x_1 \ldots x_n \mid \tau_1, \tau_2, \mu_1, \mu_2, \sigma_1^2, \sigma_2^2 \ldots)$$

$$= \prod_{i=1}^{n} \sum_{j=1}^{2} \tau_j \, f_j(x_i \mid \theta)$$

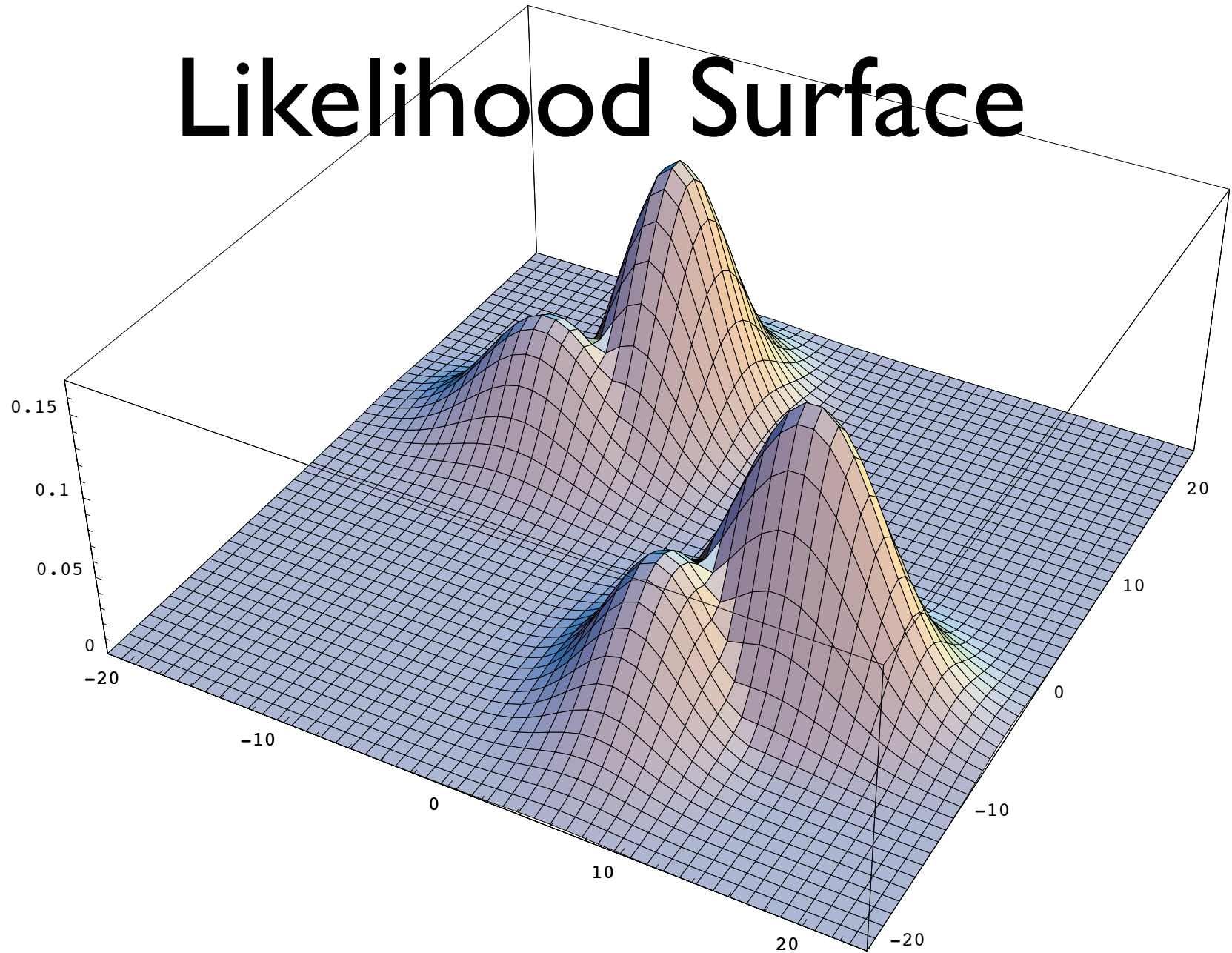Probably too messy for closed-form solution
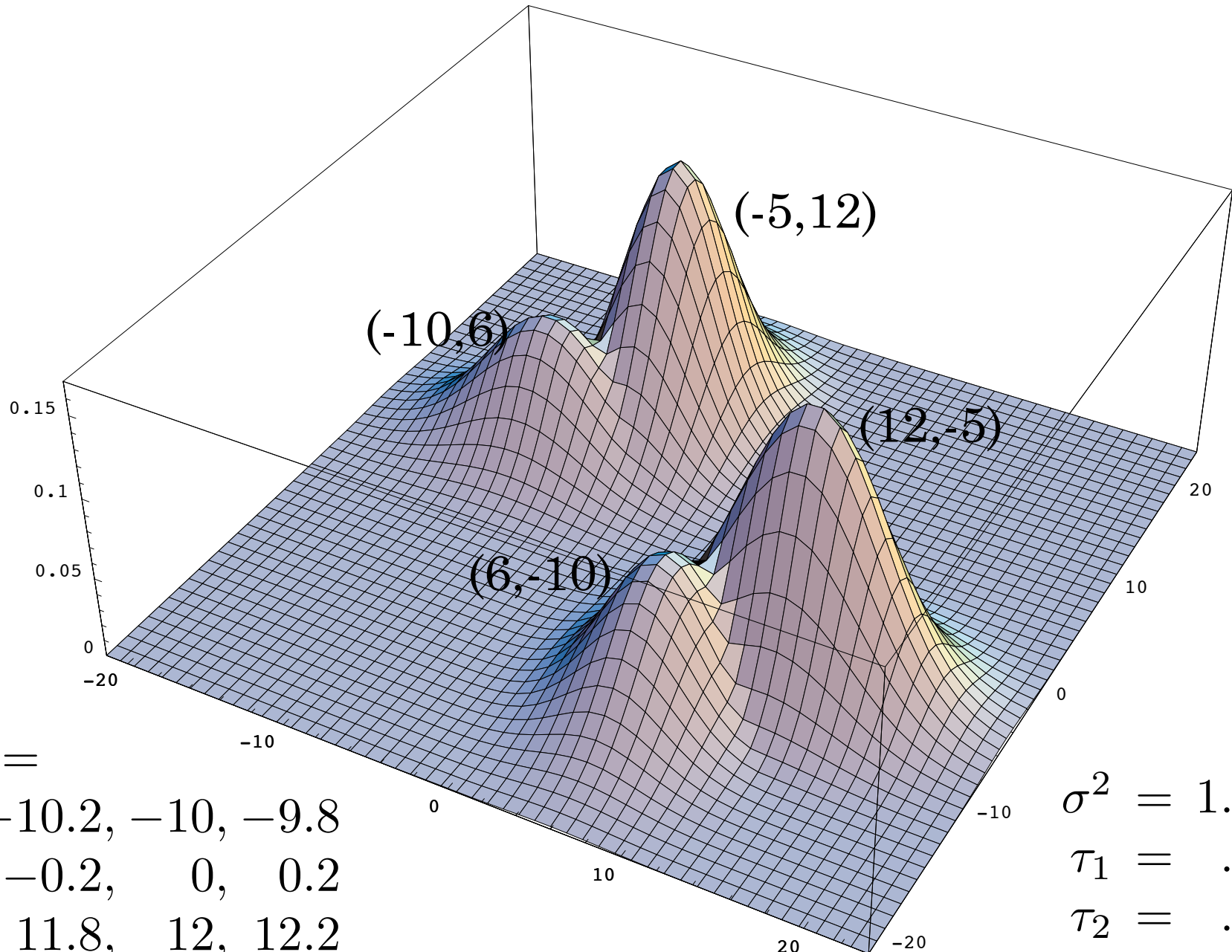
Full data

$$x_1 \quad z_{11} \quad z_{12}$$
$$x_2 \quad z_{21} \quad z_{22}$$
$$x_3 \quad z_{31} \quad z_{32}$$

$$z_{ij} \begin{cases} 0 \\ 1 \end{cases} \text{if } x_i \text{ comes from distribution } j$$

Hidden Variables

Likelihood Surface

$(-5,12)$

$(-10,6)$

$(12,-5)$

$(6,-10)$

$x_i =$
$$-10.2, \ -10, \ -9.8$$
$$-0.2, \quad 0, \quad 0.2$$
$$11.8, \quad 12, \ 12.2$$

$\sigma^2 = 1.0$
$\tau_1 = \ .5$
$\tau_2 = \ .5$

assume $\pi_g$ $\theta_j$ fixed

A event that $x_i$ drawn from $f_1$

B $\cdots$ $f_2$

D "data" $x_i$ is observed

$P(A|D)$

$P(D|A)$ Bayes rule

$P(A|D) = \dfrac{P(D|A)\,P(A)}{P(D)}$

$P(D) = P(D|A)\cdot P(A) + P(D|B)\cdot P(B)$

$f_1(x_1|\theta_1)$ $\pi_1$  $f_2(x_1|\theta_2)$ $\pi_2$

Expected value of $z_{i1}$

# M step

$$L(x_1 \; z_{11} \; z_{12} \; x_2 \; z_{21} \; z_{22} \cdots | \theta \tau)$$

$x_i$'s known

$\underline{\text{if }} z_{ij}$ know, then MLE $\theta \tau$ easy

But we don't.

Instead maximize <u>expected</u>
likelihood of visible data

$$E(L(x_1 \; x_2, \cdots x_n | \theta \tau))$$

where Expectation is over distribution
of hidden values ($z_{ij}$'s)

$$L(\vec{X}, \vec{Z} | \Theta T)$$

$$\prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} \sum_{j=1}^{2} Z_{ij} \cdot (X_i - \mu_j)^2}$$

$$E(\ln L(\vec{X}\vec{Z}|\Theta T)) =$$

$$E\left[ \sum_{i=1}^{n} -\frac{1}{2} \ln 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum_{j=1}^{2} Z_{ij} \cdot (X_i - \mu_j)^2 \right]$$

$$= \sum_{i=1}^{n} -\frac{1}{2} \ln 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum_{j=1}^{2} E(Z_{ij}) (X_i - \mu_j)^2$$

Find $\mu_j$ maximizing $\uparrow$ using $E(Z_{ij})$ from E-step.