

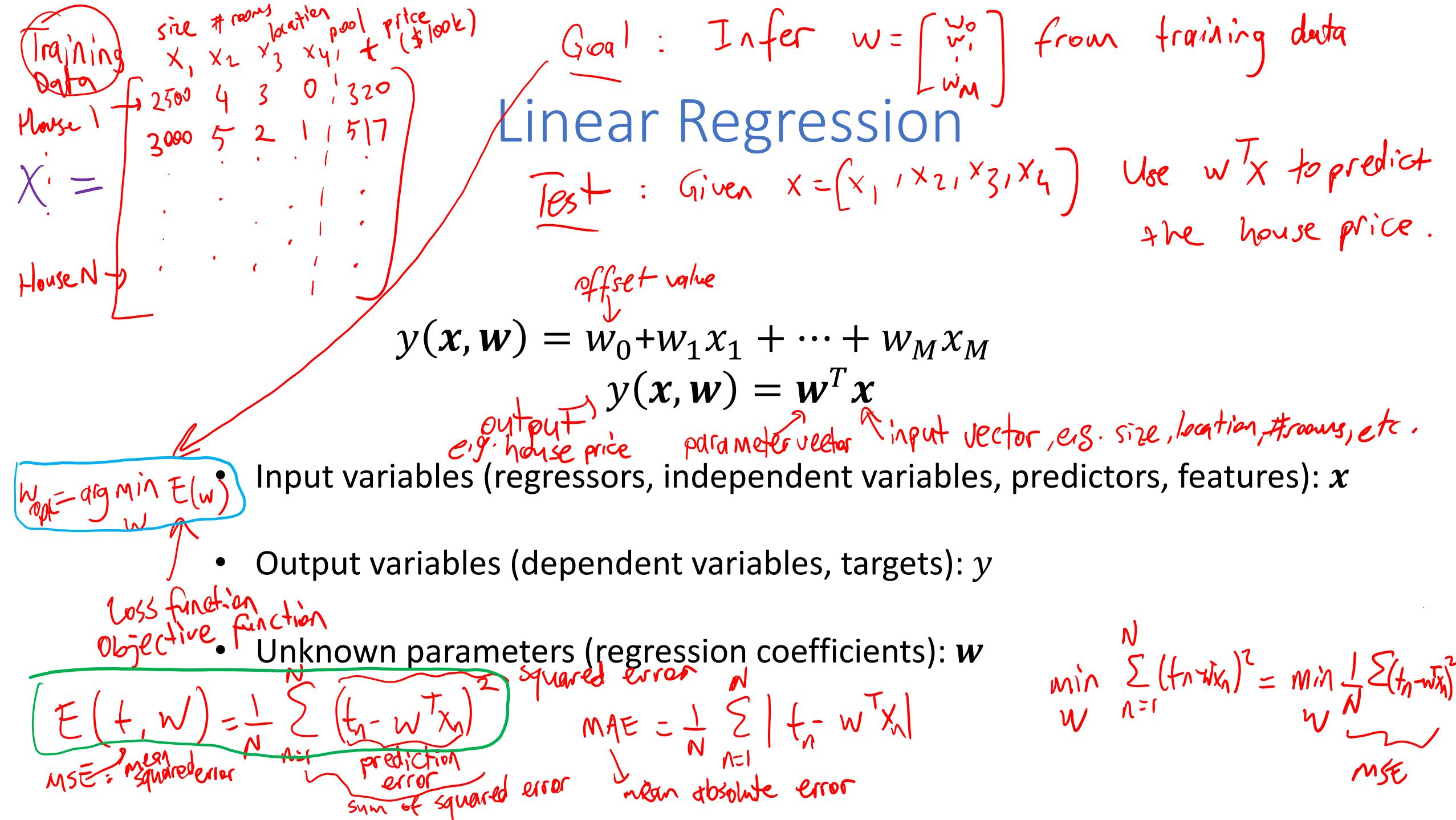
# Data Analytics

## EEE 4774 & 6777

Module 5 - Regression

Linear Regression

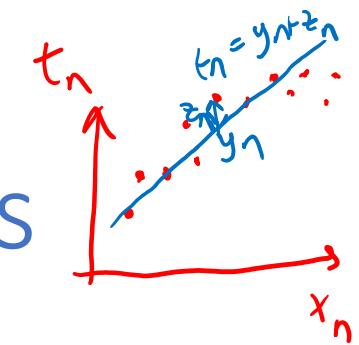
Spring 2022



# Maximum Likelihood and Least Squares

max likeli. func.

min SSE loss func.



- Assume observations from a deterministic function with added Gaussian noise:

$$t_n = \mathbf{w}^T \mathbf{x}_n + z_n \quad \text{where} \quad z_n \sim \mathcal{N}(0, \beta^{-1})$$

$\beta$ : inverse variance =  $\frac{1}{\sigma^2}$   
(precision)

which is the same as saying,

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(y(\mathbf{x}, \mathbf{w}), \beta^{-1}) \sim t_n$$

$\underbrace{\hspace{10em}}$

- Given observed inputs,  $\mathbf{X}^{N \times M}$ , and targets,  $\mathbf{t} = [t_1, \dots, t_N]^T$  we obtain the likelihood function

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \mathbf{x}_n, \beta^{-1})$$

$$X = \begin{bmatrix} -x_1 & - \\ \vdots & \vdots \\ -x_N & - \end{bmatrix}$$

# Maximum Likelihood and Least Squares

$$t_n = w^T x_n + z_n$$

$$\log p(t|w, \beta) = \sum_{n=1}^N \log \mathcal{N}(t_n | w^T x_n, \beta^{-1})$$

$$w_{ML} = \arg \min_w \frac{1}{2} \sum_{n=1}^N (t_n - w^T x_n)^2$$

$$w_{ML} = \left( \sum_{n=1}^N x_n x_n^T \right)^{-1} \sum_{n=1}^N t_n x_n$$

$$w_{ML} = (X^T X)^{-1} X^T t$$

input data matrix      output variable vector

Typically, standard Gaussian with  $\beta = 1$   
 $z_n \sim \mathcal{N}(0, 1)$

Computing the gradient and setting it to zero yields

$$w_{ML} = \frac{\sum_i x_i t_i}{\sum_i x_i^2}$$

1-dim.

$$\frac{\beta}{2\pi} e^{-\frac{\beta}{2} (t_n - w^T x_n)^2}$$

$$\min \sum_{n=1}^N (t_n - w^T x_n)^2$$

$$w_{ML} = \arg \max_w \log P(t|w, \beta)$$

ML

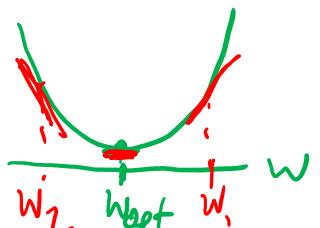
$$\sum_n x_n (x_n^T w) = \sum_n t_n x_n$$

$$(\sum_n x_n x_n^T) w = \sum_n t_n x_n$$

$$(X^T X) w = X^T t$$

$$(X^T X)^{-1} (X^T X) w = (X^T X)^{-1} X^T t$$

Sum-of-squares error  $E_D(w)$



Gradient = 0

$$\sum_{n=1}^N (t_n - w_{ML}^T x_n)^2 = \frac{1}{2} \sum_{n=1}^N (t_n - w_{ML}^T x_n)(-x_n) = 0$$

ML=LS for the Gaussian case

# Geometry of Least Squares

$$y_n = w^T \phi_n$$

- Consider

$$\mathbf{y} = X\mathbf{w}_{ML} = [\phi_1, \dots, \phi_M]\mathbf{w}_{ML}$$

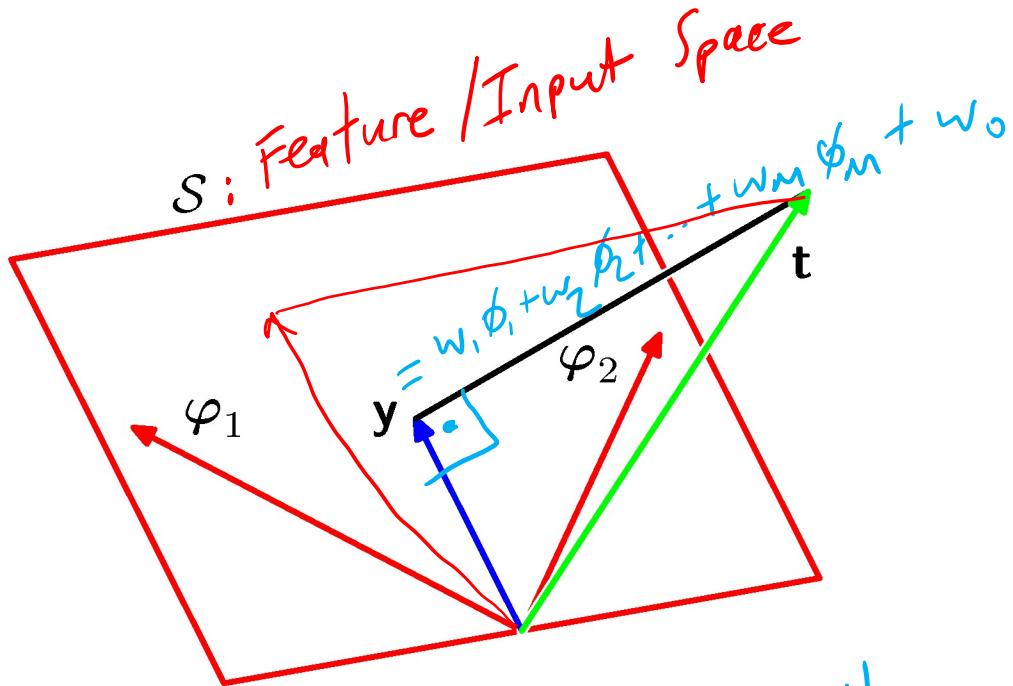
$$\vec{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

$$\mathbf{y} \in \mathcal{S} \subseteq \mathcal{T} \quad \mathbf{t} \in \mathcal{T}$$

↑ N-dimensional  
M-dimensional

- $\mathcal{S}$  is spanned by  $\phi_1, \dots, \phi_M$ .

- $\mathbf{w}_{ML}$  minimizes the distance between  $\mathbf{t}$  and its orthogonal projection on  $\mathcal{S}$ , i.e.  $\mathbf{y}$ .

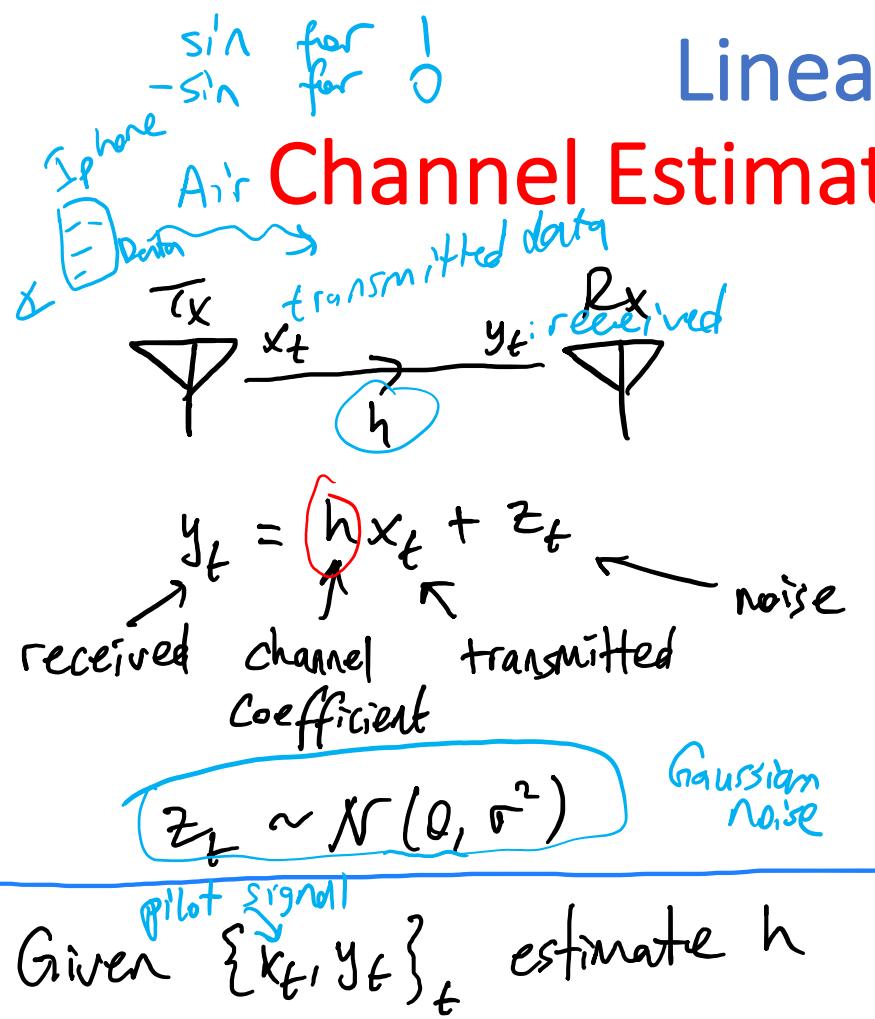


$$\|\vec{t} - \vec{y}\|^2 = \|\vec{t} - X\vec{w}\|^2 = \sum_{n=1}^N (t_n - w^T \phi_n)^2$$

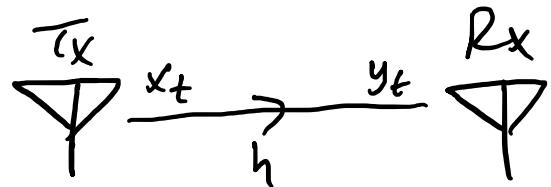
SSE

# Linear Regression Example:

## Channel Estimation in Wireless Communications



# Linear Regression Example: Channel Estimation in Wireless Communications



$$y_t = h x_t + z_t$$

received channel transmitted  
Coeficient noise

$$z_t \sim N(0, \sigma^2)$$

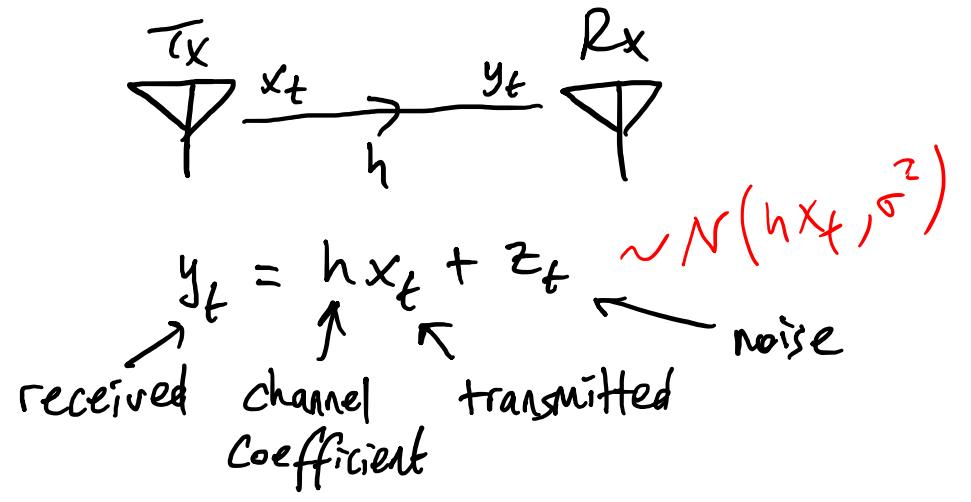
Given  $\{x_t, y_t\}_t$  estimate  $h$

Least Squares Estimation

$$\hat{h}_{LS} = \arg \min_h \sum_{t=1}^T (y_t - h x_t)^2$$
$$\frac{\partial}{\partial h} \sum_{t=1}^T (y_t - h x_t)^2 \Big|_{h=\hat{h}_{LS}} = \sum_{t=1}^T 2(y_t - \hat{h}_{LS} x_t)(-x_t) = 0 \Rightarrow$$

$$\hat{h}_{LS} = \frac{\sum_{t=1}^T y_t x_t}{\sum_{t=1}^T x_t^2}$$

# Linear Regression Example: Channel Estimation in Wireless Communications



$$z_t \sim N(0, \sigma^2)$$

Given  $\{x_t, y_t\}_t$  estimate  $h$

## Least Squares Estimation

$$\hat{h}_{LS} = \arg \min_h \sum_{t=1}^T (y_t - h x_t)^2$$

$$\frac{\partial}{\partial h} \sum_{t=1}^T (y_t - h x_t)^2 \Big|_{h=\hat{h}_{LS}} = \sum_{t=1}^T 2(y_t - \hat{h}_{LS} x_t) (-x_t) = 0 \Rightarrow$$

$$\text{Likelihood function : } p(\{x_t, y_t\} | h) = \prod_{t=1}^T p(x_t, y_t | h) = \prod_{t=1}^T N(h x_t, \sigma^2)$$

$$\boxed{\hat{h}_{LS} = \frac{\sum_{t=1}^T y_t x_t}{\sum_{t=1}^T x_t^2}}$$

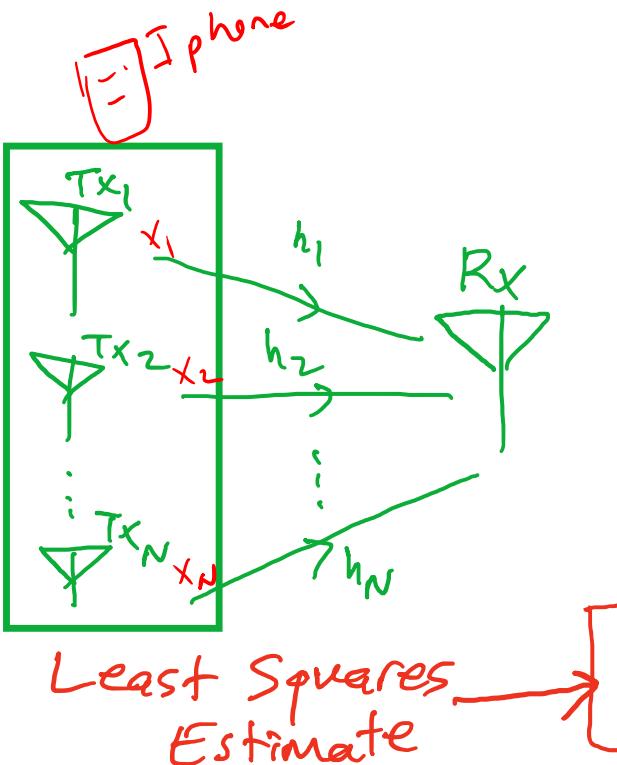
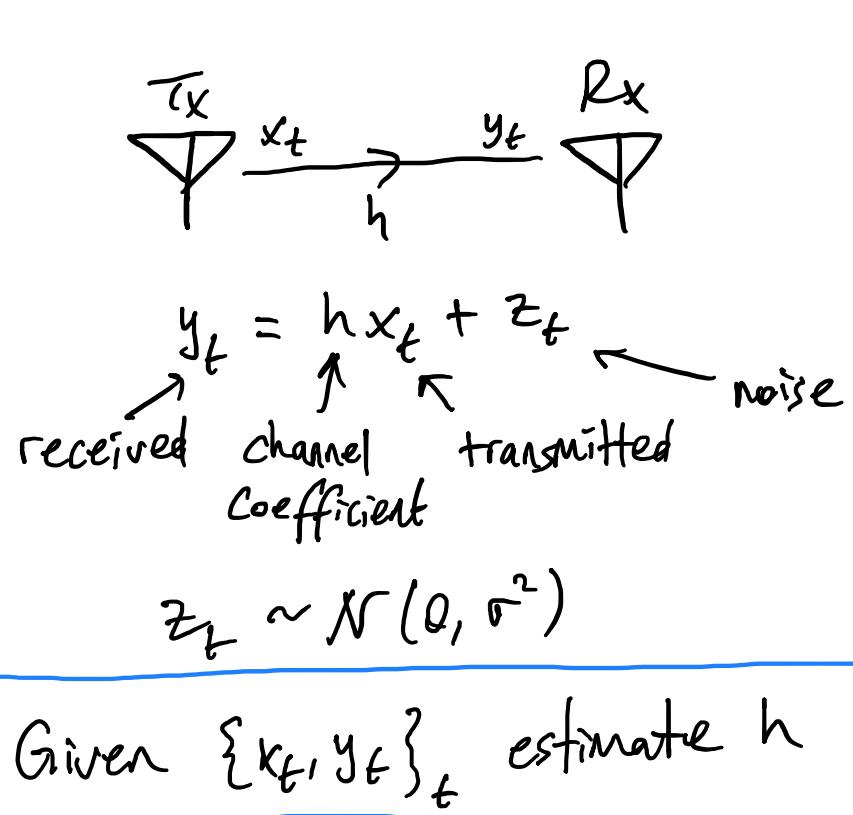
$$\prod_{t=1}^T \frac{e^{-\frac{(y_t - h x_t)^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}} = e^{-\frac{\sum_{t=1}^T (y_t - h x_t)^2}{2\sigma^2}} / (2\pi\sigma^2)^{T/2}$$

$$\text{ML Estimation} \quad \hat{h}_{ML} = \arg \max_h \log \frac{\prod_{t=1}^T (y_t - h x_t)^2}{2\pi\sigma^2}^{T/2}$$

$$\hat{h}_{ML} = \arg \min_h \sum_{t=1}^T (y_t - h x_t)^2$$

$$\boxed{\hat{h}_{ML} = \hat{h}_{LS}}$$

# Linear Regression Example: Channel Estimation in Wireless Communications



## Least Squares Estimation

$$\hat{h}_{LS} = \arg \min_h \sum_{t=1}^T (y_t - h x_t)^2$$

$$\frac{\partial}{\partial h} \sum_{t=1}^T (y_t - h x_t)^2 \Big|_{h=\hat{h}_{LS}} = \sum_{t=1}^T 2(y_t - \hat{h}_{LS} x_t)(-x_t) = 0 \Rightarrow \hat{h}_{LS} = \frac{\sum_{t=1}^T y_t x_t}{\sum_{t=1}^T x_t^2}$$

$$\hat{h}_{LS} = \frac{\sum_{t=1}^T y_t x_t}{\sum_{t=1}^T x_t^2}$$

$$\hat{h}_{ML} = \arg \max_h \log \hat{P} \frac{\sum_{t=1}^T (y_t - h x_t)^2}{2\pi\sigma^2}$$

$$\hat{h}_{ML} = \arg \min_h \sum_{t=1}^T (y_t - h x_t)^2$$

$$\hat{h}_{ML} = \hat{h}_{LS}$$

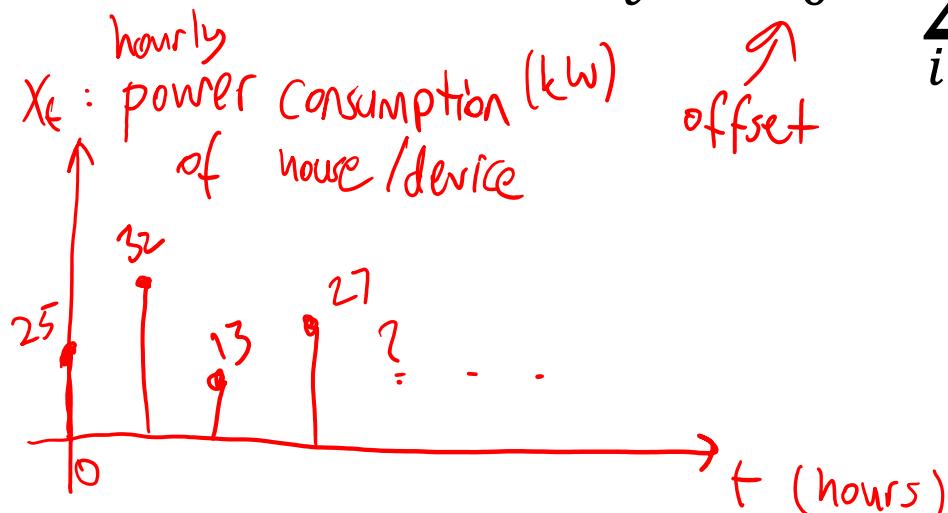
# Linear Regression Example: Autoregressive (AR) Model

- AR(p): the next value depends on the previous p values

$$x_t = w_0 + w_1 x_{t-1} + w_2 x_{t-2} + \cdots + w_p x_{t-p} + z_t$$

$$x_t = w_0 + \sum_{i=1}^p w_i x_{t-i} + z_t$$

↑ error  
↑ offset  
↑ parameters of lin. reg. model



# Regularized Least Squares

- Consider the error function:

$$\min \left\{ \beta E_D(\mathbf{w}) + \lambda E_W(\mathbf{w}) \right\}$$

Data term + Regularization term

SSE (MSE), Likelihood

- With the sum-of-squares error function and a quadratic regularizer, we get

$$\min \left\{ \frac{\beta}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \right\}$$

inverse variance in ML approach

- which is minimized by

$$\mathbf{w}_{LS} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

w/o regularization

$$\mathbf{w} = \left( \frac{\lambda \mathbf{I}}{\beta} + \Phi^T \Phi \right)^{-1} \Phi^T \mathbf{t}$$

$$\beta = 1 \quad \mathbf{w} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

$\lambda$  is called the regularization coefficient.

$\min_{\mathbf{w}} E_D(\mathbf{w})$

Choose  $\mathbf{w}$  values to minimize loss function

Complex models with many parameters, have a very high-dimensional search space for optimizing parameters. In turn, training loss can be very well minimized (i.e., over fitting to train. data happens)

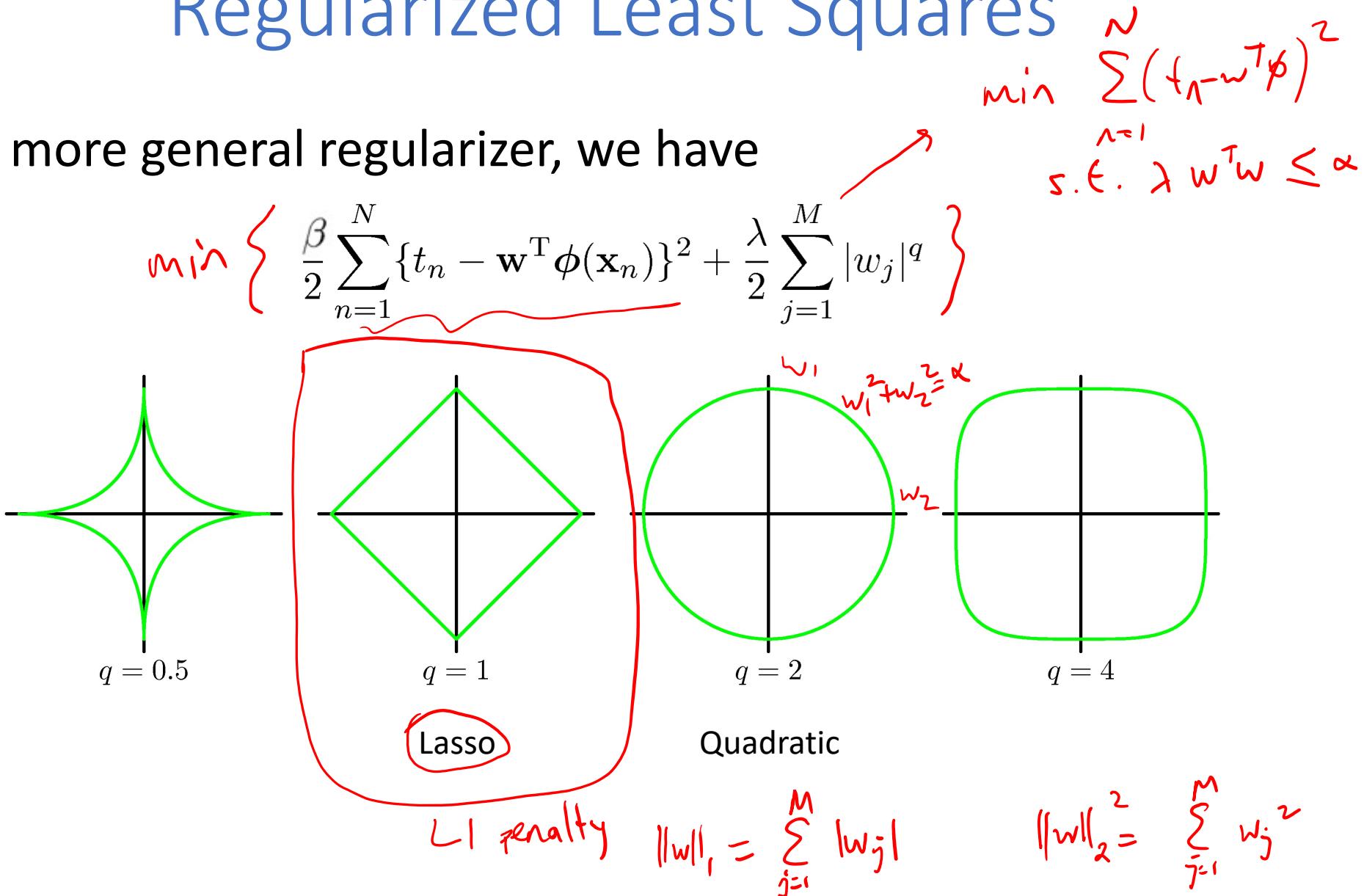
in this high-dim. unrestricted search space.

by choosing extreme values for parameters.

$$\text{e.g. } \mathbf{w} = \begin{bmatrix} 10^{12} \\ -10^{15} \\ \vdots \end{bmatrix}$$

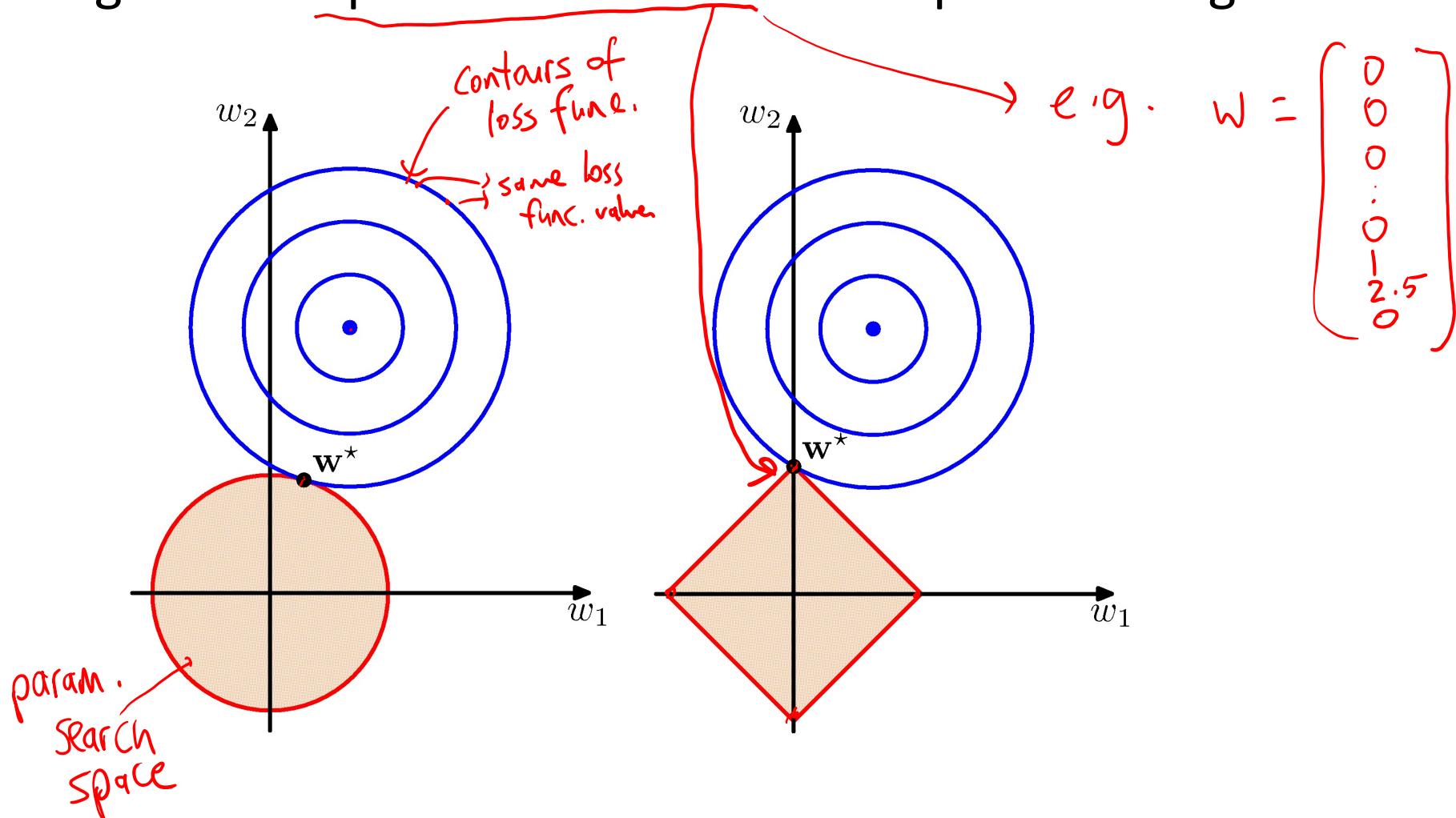
# Regularized Least Squares

- With a more general regularizer, we have



# Regularized Least Squares

- Lasso tends to generate sparser solutions than a quadratic regularizer.



$$MSE = E\{y(x) - t\}^2 = E\{y(x) - h(x) + h(x) - t\}^2 = E\{y(x) - h(x)\}^2 + E\{h(x) - t\}^2 + 2E\{y(x) - h(x)\} \cancel{E\{h(x) - t\}}$$

How to select regularization coefficient  $\lambda$ ?

## The Bias-Variance Decomposition

- Recall the *expected squared loss*, i.e., *Mean Squared Loss (MSE)*,

$$MSE = \mathbb{E}[L] = \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \underbrace{\iint \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt}_{\text{given information}}$$

where

$$h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] = \int tp(t|\mathbf{x}) dt.$$

Best Bayesian predictor  
 that minimizes MSE  
 given information  
 posterior mean  
 posterior prob. distn. of output var. given input  
 var.

$$\begin{aligned} & \text{Cross-term} = 0 \\ & E[y(\mathbf{x}) - h(\mathbf{x})] E[h(\mathbf{x}) - t] \end{aligned}$$

$$E[h(\mathbf{x})] = E_x[\mathbb{E}[t|\mathbf{x}]]$$

$$\text{Tower property} = E[t]$$

$$\mathbb{E}[h(\mathbf{x}) - t] = \mathbb{E}[t] - E[t] = 0$$

- The second term of  $E[L]$  corresponds to the noise inherent in the random variable  $t$ .
- What about the first term?

$$MSE = E\{y(x) - h(x)\}^2 + E\{h(x) - t\}^2$$

Focus on this

# The Bias-Variance Decomposition

- Suppose we were given multiple data sets, each of size N. Any particular data set, D, will give a particular function  $y(x; D)$ . We then have

$$\begin{aligned} & \mathbb{E} \left[ \{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2 \right] \\ &= \mathbb{E} \left[ \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] + \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \right] \\ &= \mathbb{E} \left[ \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2 + \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \right. \\ &\quad \left. + 2\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}. \right] \end{aligned}$$

*Variance of our predictor  $y$*

*Bias of our predictor  $y$*

$$= \mathbb{E} \left[ \{y - \underbrace{\mathbb{E}[y]}_{\text{mean}}\}^2 \right] + \mathbb{E} \left[ \{ \underbrace{\mathbb{E}[y] - h(\mathbf{x})}_{\text{Bias}} \}^2 \right]$$

# The Bias-Variance Decomposition

- Taking the expectation over  $\mathcal{D}$  yields

*Controllable part of MSE*  $\rightarrow \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2]$

$$= \underbrace{\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2}_{\text{(bias)}^2} + \underbrace{\mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2]}_{\text{variance}}.$$

# The Bias-Variance Decomposition

- Thus we can write

$$\text{expected loss} = (\text{bias})^2 + \text{variance} + \text{noise}$$

where

*depend on predictor*

$$\begin{array}{c} \min \\ \text{Two tasks} \end{array} \quad \begin{array}{c} \min \\ \cdot \end{array}$$

$$(\text{bias})^2 = \int \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x}$$

$$\text{variance} = \int \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2] p(\mathbf{x}) d\mathbf{x}$$

$$\text{noise} = \iint \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

$$y(\tilde{x}) = w^T \tilde{x}$$

$$\tilde{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

$$y(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3 + \dots + w_m x^m$$

M=1 :  $y(x) = w_0 + w_1 x$

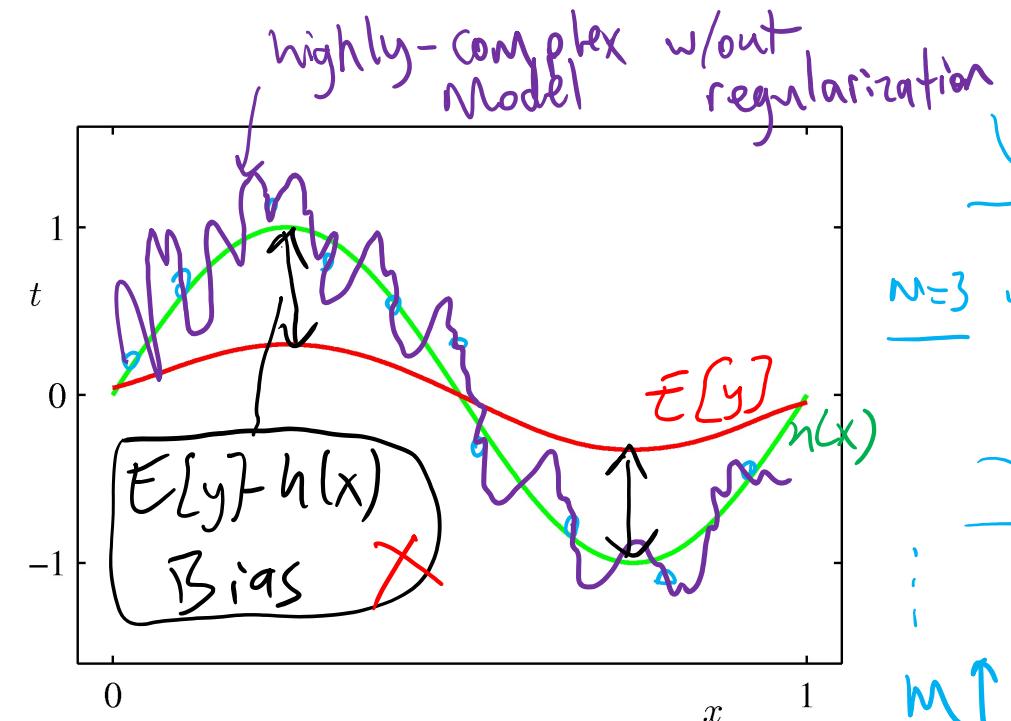
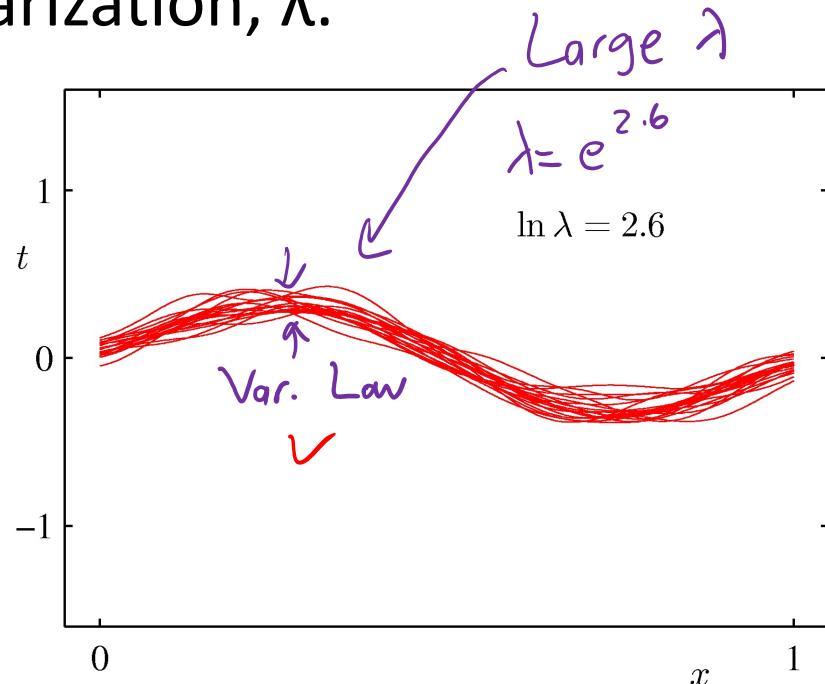


## The Bias-Variance Decomposition

$$\min \sum_{i=1}^n (y(x_i) - f_i)^2 + \frac{\lambda}{2} \|w\|^2$$

M=2 :  $y(x) = w_0 + w_1 x + w_2 x^2$

- Example: 25 data sets from the sinusoidal, varying the degree of regularization,  $\lambda$ .



M=3 :  $y(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3$

M=4 :  $y(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4$

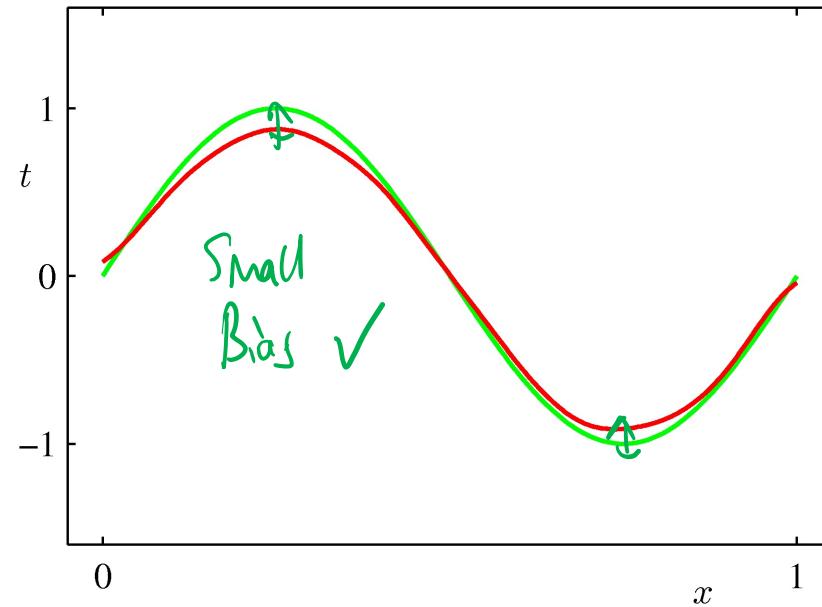
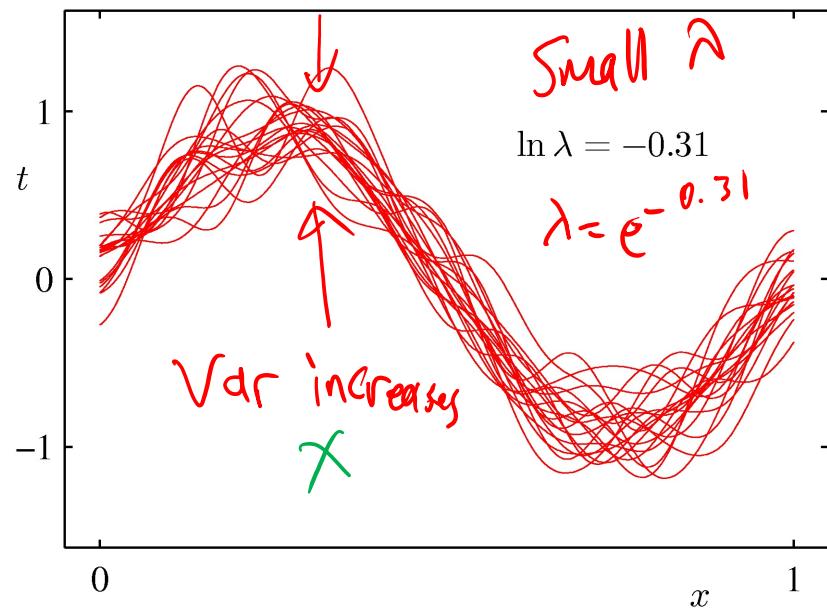
$M \uparrow$  ↑ modeling power (model complexity) flexibility

Note :

Complex models w/ many parameters should only be used for complex problems w/ abundance of data  
 to estimate in training

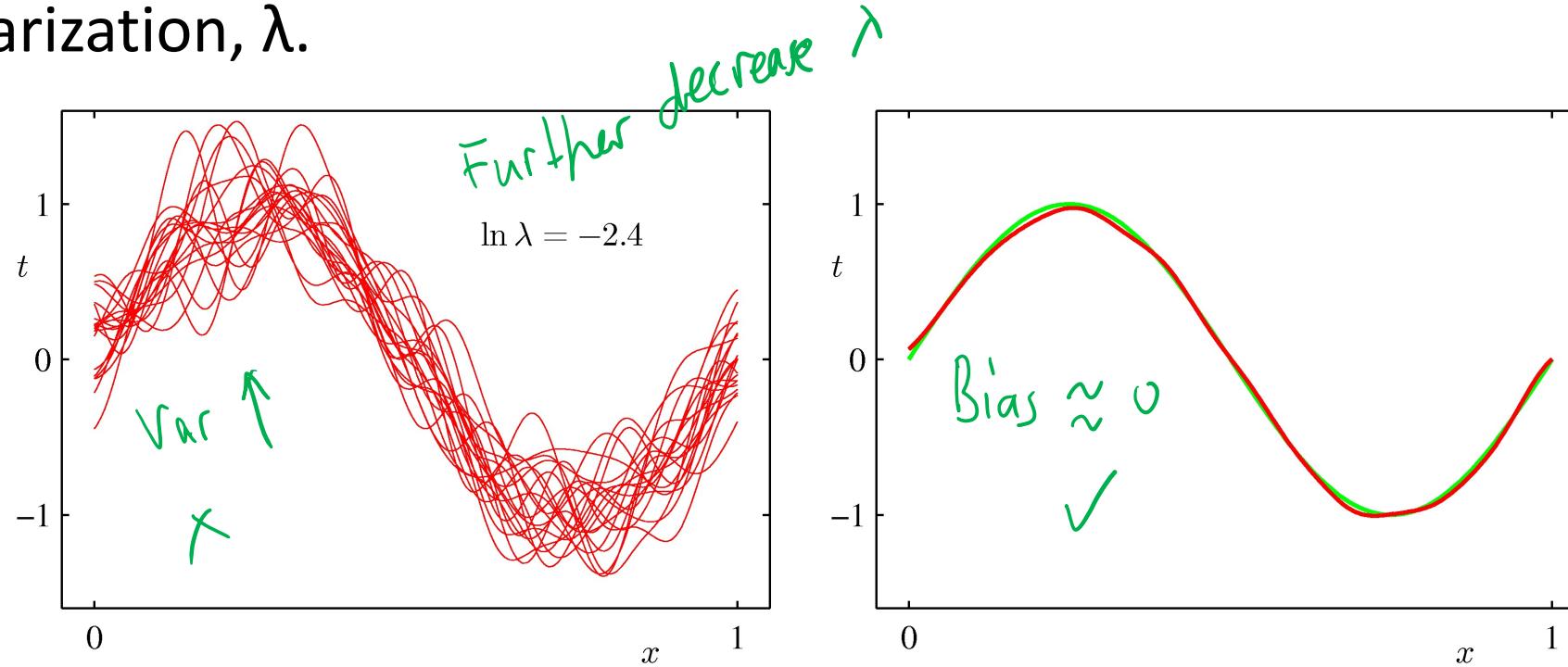
# The Bias-Variance Decomposition

- Example: 25 data sets from the sinusoidal, varying the degree of regularization,  $\lambda$ .



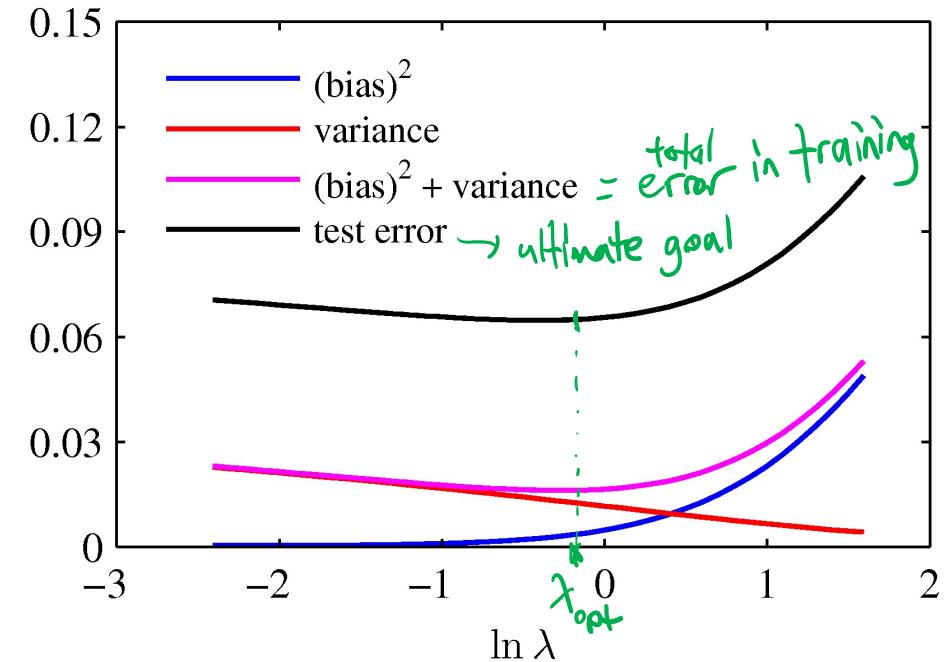
# The Bias-Variance Decomposition

- Example: 25 data sets from the sinusoidal, varying the degree of regularization,  $\lambda$ .



# The Bias-Variance Trade-off

- From these plots, we note that an over-regularized model (large  $\lambda$ ) will have a high bias, while an under-regularized model (small  $\lambda$ ) will have a high variance.



$$\mathbf{t} = \mathbf{w}^\top \mathbf{x} + \mathbf{z}$$

# Bayesian Linear Regression

- Define a conjugate prior over  $\mathbf{w}$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0).$$

prior mean

prior covar.

- Combining this with the likelihood function and using results for marginal and conditional Gaussian distributions, gives the posterior

- where

$$p(\mathbf{w} | \mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)$$

$$\begin{aligned}\mathbf{m}_N &= \mathbf{S}_N \left( \mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t} \right) \\ \mathbf{S}_N^{-1} &= \mathbf{S}_0^{-1} + \beta \Phi^T \Phi.\end{aligned}$$

posterior

data matrix

# Bayesian Linear Regression

- A common choice for the prior is

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I})$$
$$\mathbf{S}_0 = \begin{bmatrix} \alpha^{-1} & & & \\ & \alpha^{-1} & & \\ & & \ddots & \\ & & & \alpha^{-1} \end{bmatrix}$$

- for which

$$\mathbf{m}_N = \beta \underline{\mathbf{S}_N} \Phi^T \mathbf{t} = \beta (\alpha \mathbf{I} + \beta \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$
$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi.$$
$$= (\underbrace{\alpha \mathbf{I} + \beta \Phi^T \Phi}_{\frac{\beta}{\alpha}})^{-1} \Phi^T \mathbf{t}$$

- Next we consider an example ...

$$\hat{\mathbf{w}}_{\text{Bay}} = \mathbf{m}_N$$
$$= w_{LS + \text{Quad Reg}}$$

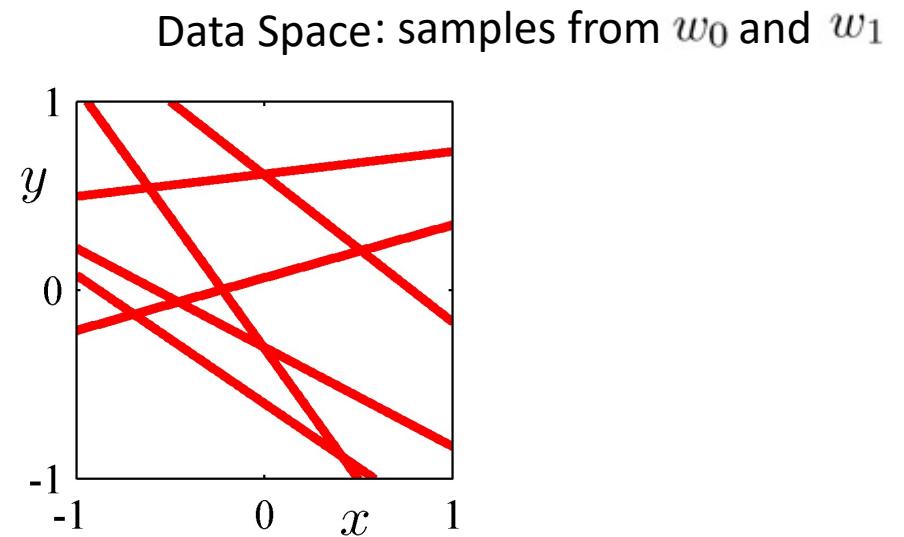
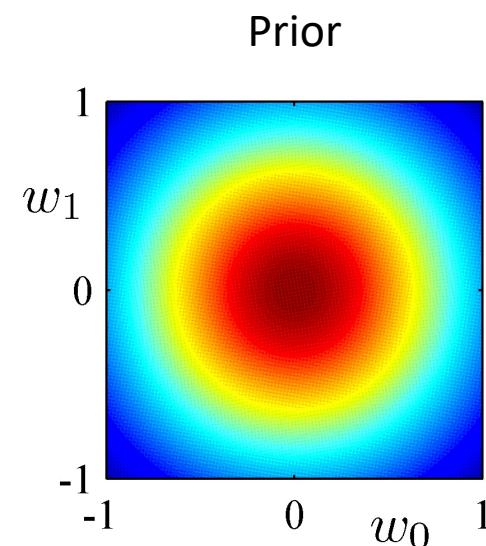
$$\alpha \equiv \lambda$$

↑  
reg. coef.

# Bayesian Linear Regression

0 data points observed

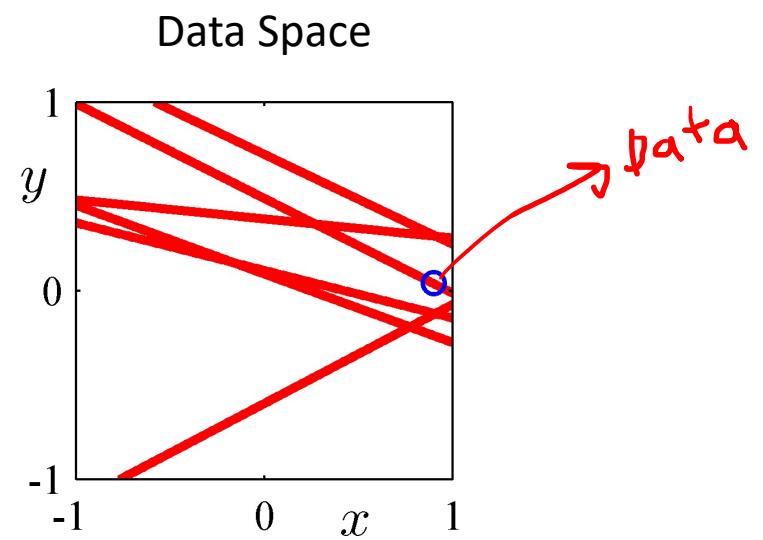
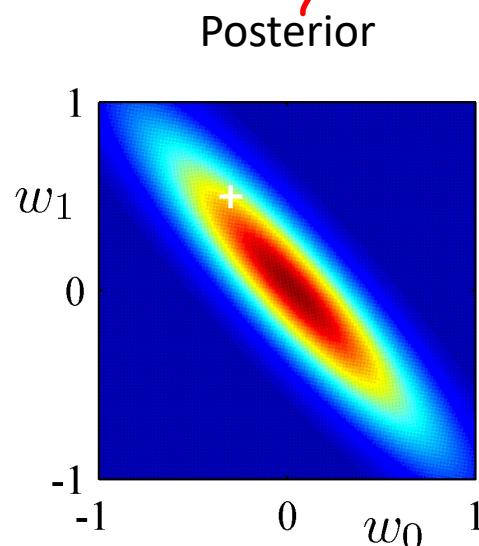
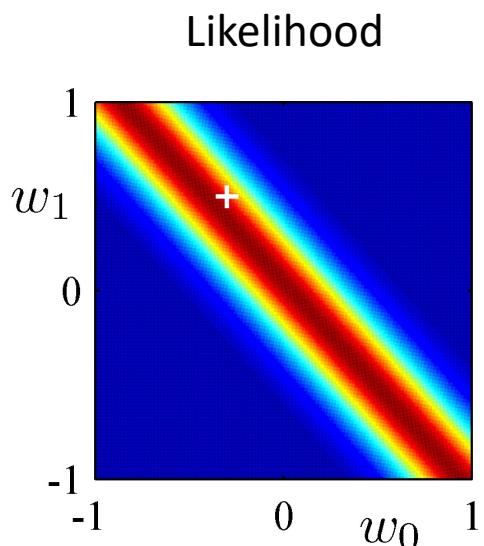
$$y = w_0 + w_1 x$$



$$\mathcal{S}_0 = \begin{bmatrix} \alpha^{-1} & 0 \\ 0 & \alpha^{-1} \end{bmatrix}$$

# Bayesian Linear Regression

1 data point observed



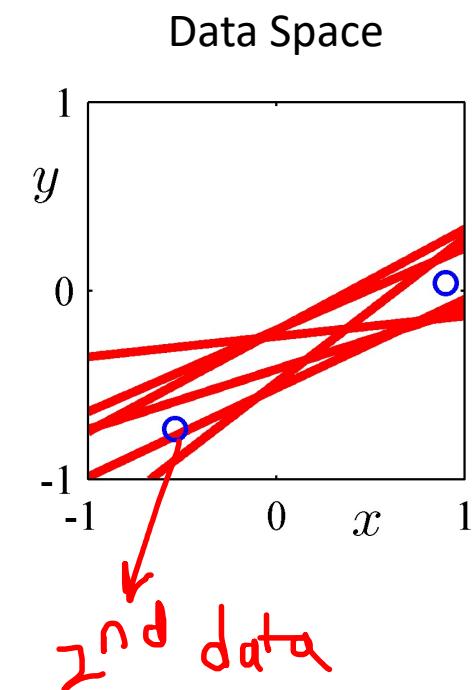
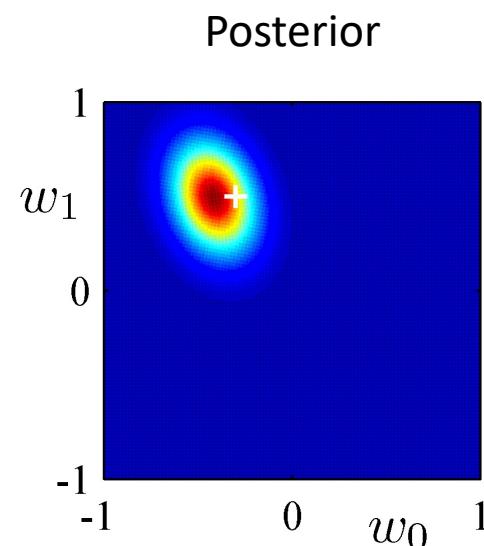
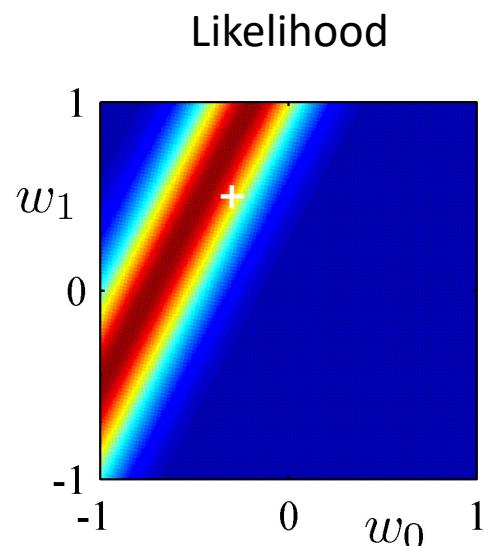
Likelihood . Prior

prior for 2<sup>nd</sup> data

data

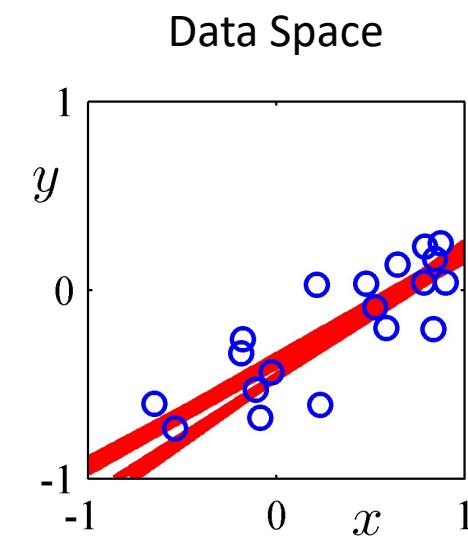
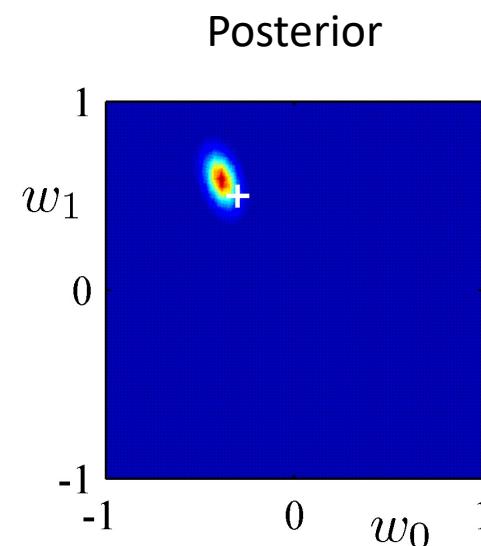
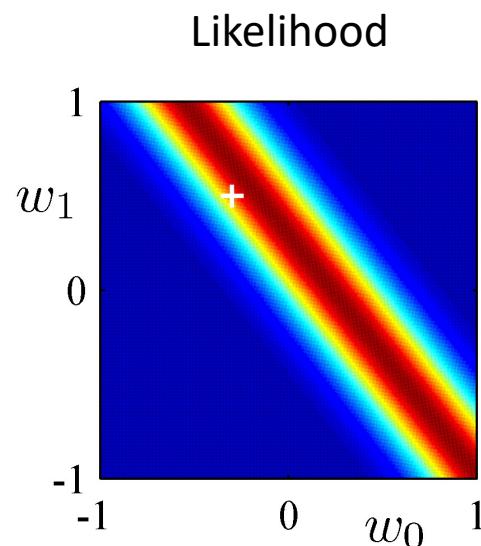
# Bayesian Linear Regression

2 data points observed



# Bayesian Linear Regression

20 data points observed



# Sequential Learning

$$MSE = \mathbb{E}[(t_n - w^T \phi(x_n))^2], \quad \hat{w}_{MSE} = \arg \min_w MSE(w)$$

$$\Rightarrow \hat{w}_{MSE} = \mathbb{E}[\phi \phi^T]^{-1} \mathbb{E}[t \phi]$$

$$\begin{aligned}\hat{w}_{LS} &= \arg \min_w \sum_{n=1}^N (t_n - w^T \phi(x_n))^2 \\ &= (\sum_n t t^T)^{-1} \sum_n \phi \phi^T\end{aligned}$$

Batch Learning

- Data items considered one at a time (a.k.a. online learning); use stochastic (sequential) gradient descent:

Gradient Descent: Iterative optimization – move towards negative of gradient of cost function MSE

Stochastic Gradient Descent: Use single-point estimate  $(t_n - w^T \phi) \phi$

$$\begin{aligned}w^{(\tau+1)} &= w^{(\tau)} - \eta \nabla E_n \\ &= w^{(\tau)} + \eta(t_n - w^{(\tau)} \phi(x_n)) \phi(x_n).\end{aligned}$$

$\frac{1}{N} \sum_{n=1}^N (t_n - w^T \phi) \phi$

Can be estimated by

$$\begin{aligned}&\mathbb{E}[\nabla_w (t_n - w^T \phi)^2] \\ &= \mathbb{E}[2(t_n - w^T \phi)(-\phi)]\end{aligned}$$

Difficult to compute online

- This is known as the *least-mean-squares (LMS) algorithm*. Issue: how to choose  $\eta$ ?