

# Linear models for regression } Lecture 2

Simplest model:

$$y(\vec{x}, \vec{w}) = w_0 + \sum_{j=1}^D w_j x_j$$

$\{\vec{x}_1, \dots, \vec{x}_N\}$  <sup>multi-D,</sup> observations  $\rightarrow$  each  $\vec{x} = \{x_1, \dots, x_D\}$   
 $\{t_1, \dots, t_N\}$  target variables

More generally,

$$y(\vec{x}, \vec{w}) = \underbrace{w_0}_{\text{bias prm}} + \sum_{j=1}^D w_j \underbrace{f_j(\vec{x})}_{\text{basis f's}}$$

still a linear model because  $y$  is linear in  $\vec{w}$

Sometimes, one defines  $f_0(\vec{x}) = 1$ , s.t.

$$y(\vec{x}, \vec{w}) = \sum_{j=0}^D w_j f_j(\vec{x}) = \vec{w}^T \vec{f}(\vec{x})$$
$$\vec{w} = (w_0 \dots w_D) \quad \vec{f} = (f_0 \dots f_D)$$

Popular choices of basis f's:  
(1D examples)

(1) Gaussian  $f_j(x) = e^{-\frac{(x-\mu_j)^2}{2\sigma^2}}$

(2) Sigmoidal  $f_j(x) = \sigma\left(\frac{x-\mu_j}{s}\right)$ , where

$$\sigma(a) = \frac{1}{1+e^{-a}}$$

As before, we assume that

$$t = y(\vec{x}, \vec{w}) + \xi$$

$\uparrow$  target       $\uparrow$  model       $\uparrow$  noise

$$P(\xi | \beta) = \mathcal{N}(0, \beta^{-1})$$

Correspondingly,

$$P(t | \vec{x}, \vec{w}, \beta) = \mathcal{N}(t | y(\vec{x}, \vec{w}), \beta^{-1})$$

Now consider

$$\vec{X} = \{\underbrace{\vec{x}_1, \dots, \vec{x}_N}_{\text{inputs}}\}$$

$$\vec{T} = \{\underbrace{t_1, \dots, t_N}_{\text{target vars}}\}$$

Under the independence assumption,

$$P(\vec{T} | \vec{X}, \vec{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \vec{w}^T \cdot \vec{y}(\vec{x}_n), \beta^{-1})$$

$\downarrow$   
 drop conditional dependence  
 on  $\vec{X}$  for brevity

Then

$$\log P(\vec{T} | \vec{w}, \beta) = \frac{N}{2} \log \beta - \frac{N}{2} \log(2\pi) - \beta \underbrace{\frac{1}{2} \sum_{n=1}^N (t_n - \vec{w}^T \vec{y}(\vec{x}_n))^2}_{E_{\bullet}(\vec{w})} \quad (**)$$

$$\stackrel{N}{\sum_{n=1}^N} \log \mathcal{N}(t_n | \vec{w}^T \vec{y}(\vec{x}_n), \beta^{-1})$$

Now, use ML to determine  $\vec{w}_{ML}$  &  $\beta_{ML}$ .

Start with the weights:

$$\frac{\partial}{\partial w_j} \log P(\vec{t} | \vec{w}, \beta) = \frac{\beta}{2} \sum_{n=1}^N (t_n - \vec{w}^T \vec{\psi}(\vec{x}_n)) \times \psi_j(\vec{x}_n),$$

or

$$\underbrace{\vec{\nabla}_{\vec{w}}}_{=0} \log P = \beta \sum_{n=1}^N (t_n - \vec{w}^T \vec{\psi}(\vec{x}_n)) \vec{\psi}(\vec{x}_n) = 0$$

$$\Rightarrow \sum_{n=1}^N t_n \psi_j(\vec{x}_n) = \sum_{i=0}^D w_i \sum_{n=1}^N \psi_i(\vec{x}_n) \psi_j(\vec{x}_n)$$

Call  $\psi_j(\vec{x}_n) = \phi_{nj} \rightarrow$  elements of  $N \times (D+1)$  design matrix

Then  $\underbrace{\sum_n t_n \phi_{nj}}_{D+1 \text{ vector}} = \sum_i w_i \underbrace{\sum_n \phi_{ni} \phi_{nj}}_{(D+1) \times (D+1) \text{ matrix}} = \underbrace{\sum_i (\phi^T \phi)_{ji} w_i}_{D+1 \text{ vector}}$

Finally,  $\phi^T \vec{t} = (\phi^T \phi) \vec{w}$ , or

$$\vec{w}_{ML} = (\phi^T \phi)^{-1} \phi^T \vec{t}$$



normal eq's for the least-squares problem

$$\phi = \begin{pmatrix} \psi_0(\vec{x}_1) & \psi_1(\vec{x}_1) & \dots & \psi_D(\vec{x}_1) \\ \vdots & \vdots & \ddots & \vdots \\ \psi_0(\vec{x}_N) & \psi_1(\vec{x}_N) & \dots & \psi_D(\vec{x}_N) \end{pmatrix}$$

$\tilde{\Phi} = (\Phi^T \Phi)^{-1} \Phi^T$  is called the Moore-Penrose pseudo-inverse of  $\Phi$ .

---

Next, consider

$$E(\vec{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - w_0 - \sum_{j=1}^D w_j \phi_j(\vec{x}_n))^2$$

$$\frac{\partial E}{\partial w_0} = - \sum_{n=1}^N (t_n - w_0 - \sum_{j=1}^D w_j \phi_j(\vec{x}_n)) = 0$$

gives

$$w_0^{ML} = \frac{1}{N} \sum_{n=1}^N t_n - \frac{1}{N} \sum_{n=1}^N \sum_{j=1}^D w_j \phi_j(\vec{x}_n) =$$

$$= \bar{t} - \sum_{j=1}^D w_j \bar{\phi}_j, \quad \text{where } \bar{\quad} \text{ indicates } \frac{1}{N} \sum_{n=1}^N \dots$$

Thus  $w_0^{ML}$  is the difference between the average of the target values and the weighted sum of basis f'n averages.

---

Finally, maximize (\*\*) wrt  $\beta$ :

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N (t_n - \vec{w}_{ML}^T \vec{\phi}(\vec{x}_n))^2$$

residual variance of the target values around

$$y(\vec{x}, \vec{w})$$

# Geometry of least squares

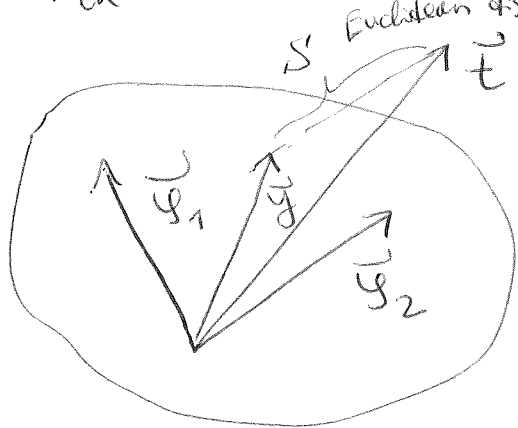
$\vec{t} = (t_1 \dots t_N)$  is an  $N$ -dim vector

each  $\vec{y}_j = (y_j(\vec{x}_1) \dots y_j(\vec{x}_N))$  is also an  $N$ -dim vector

[ $\vec{y}_j$  is the  $j^{\text{th}}$  column of  $\Phi$ ]

If  $D+1 < N$ , <sup>#datapoints</sup> then  $\vec{y}_j$  vectors the # of basis functions, including  $y_0(\vec{x}_n)$

will span a subspace  $S$  of dimensionality  $D+1$  in the  $N$ -dim'l space.



Define  $\vec{y} = (y(\vec{x}_1, \vec{w}) \dots y(\vec{x}_N, \vec{w}))$

$N$ -dim vector, a linear combination of

$D+1$   $\vec{y}_j$  vectors.

$E(\vec{w}) \sim$  ~~square of~~ square of Euclidean dist. between  $\vec{y}$  &  $\vec{t}$ , so we need to minimize it. This is accomplished by the orthogonal projection of  $\vec{t}$  into subspace  $S$ .

In practice,  $\Phi^T \Phi$  may be hard to invert, when some  $\vec{y}_j$ 's are nearly colinear  $\rightarrow$  add a regularization term to prevent this.

# Regularized least squares

Consider now

$$\tilde{E}(\vec{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \vec{w}^T \vec{\phi}(\vec{x}_n))^2 + \frac{\lambda}{2} \vec{w}^T \vec{w}$$

↑ still a quadratic f'n of  $w_j$

$$\vec{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \vec{t} \quad \text{becomes}$$

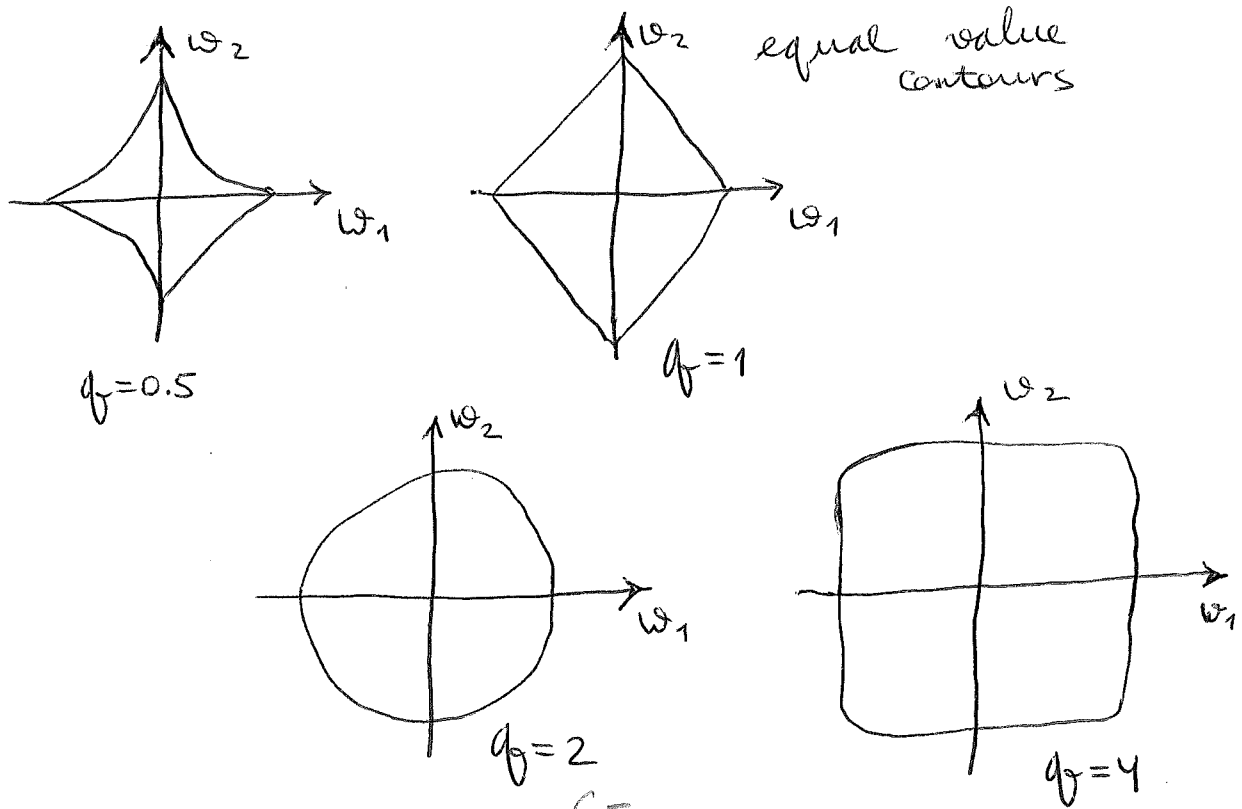
$$\vec{w}^* = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \vec{t} \quad \left[ \begin{array}{l} \uparrow \text{unit matrix } (D+1) \times (D+1) \\ \text{the inverse matrix is now regularized} \end{array} \right]$$

More generally,

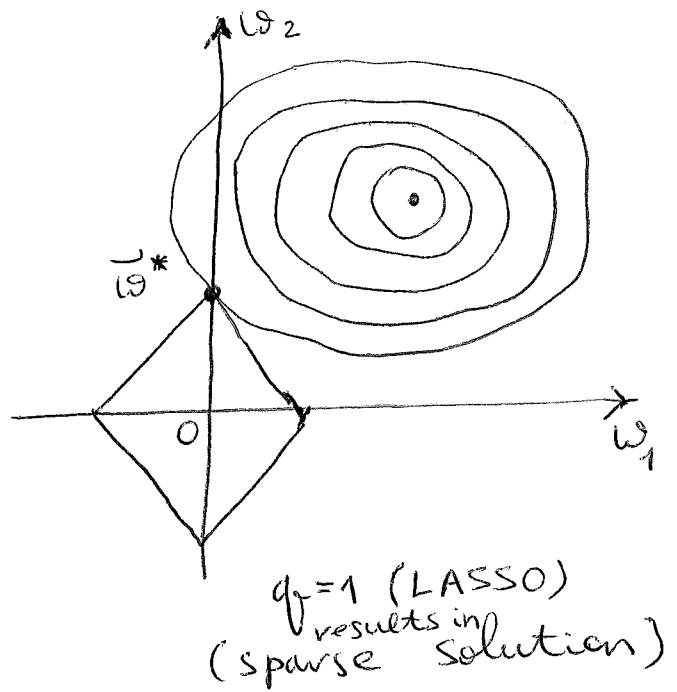
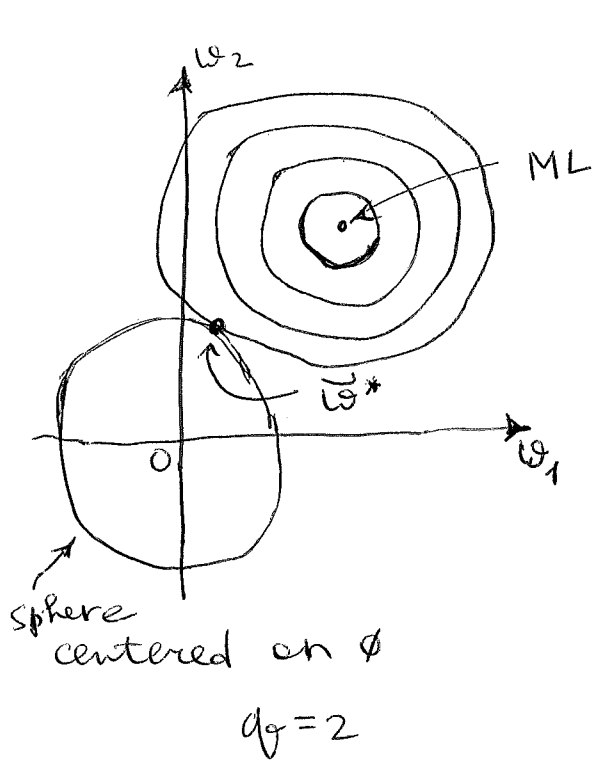
$$\tilde{E}(\vec{w}) = \frac{1}{2} \sum_{n=1}^N (\dots)^2 + \frac{\lambda}{2} \sum_{j=1}^{D+1} |w_j|^{q_f}$$

$q_f = 2 \leftarrow$  quadratic regularizer above

$q_f = 1 \leftarrow$  LASSO regularizer



Now, consider



### Loss functions for regression (1.5.5)

Consider  $\vec{x} \rightarrow t$   
input target

$$y(\vec{x}) \xrightarrow{\text{model for } t} L(t, y(\vec{x})) \stackrel{\text{e.g.}}{=} \underbrace{(y(\vec{x}) - t)^2}_{\text{loss}} \quad \text{squared loss}$$

$$E[L] = \int d\vec{x} dt \underbrace{p(\vec{x}, t)}_{\text{joint prob. of } \vec{x} \& t} \underbrace{(y(\vec{x}) - t)^2}_{\text{loss f'n}}$$

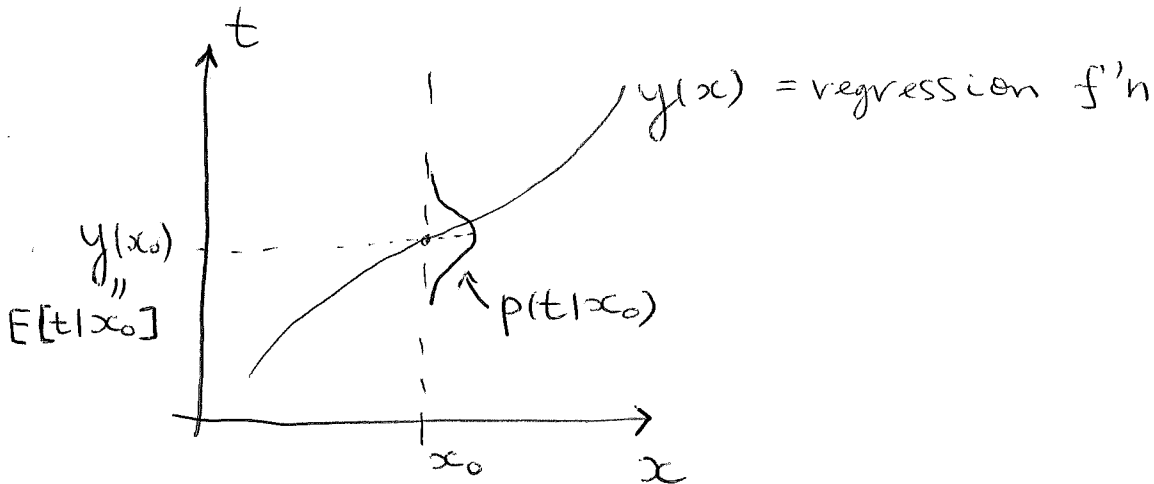
expected value of loss

$$\frac{\delta E[L]}{\delta y(\vec{x})} = 2 \int dt (y(\vec{x}) - t) p(\vec{x}, t) = 0$$

$\rightarrow$  minimize  $E[L]$

$$y(\vec{x}) = \frac{\int dt p(\vec{x}, t) t}{p(\vec{x})} = \int dt t p(t|\vec{x}) = E[t|\vec{x}] \quad (+)$$

$y(\vec{x})$  is called a regression f'n



Further,

$$(y(\vec{x}) - t)^2 = (y(\vec{x}) - E[t|\vec{x}] + E[t|\vec{x}] - t)^2 =$$

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

f'n of  $\vec{x}$ , not  $t$

Then 
$$E[L] = \int dt d\vec{x} p(\vec{x}, t) \left[ (y(\vec{x}) - E[t|\vec{x}])^2 + 2(y(\vec{x}) - E[t|\vec{x}])(E[t|\vec{x}] - t) + (E[t|\vec{x}] - t)^2 \right] =$$

$$= \int d\vec{x} p(\vec{x}) \left[ (y(\vec{x}) - E[t|\vec{x}])^2 + \int dt p(\vec{x}, t) p(\vec{x}) \right]$$

$$+ 2 \int d\vec{x} (y(\vec{x}) - E[t|\vec{x}]) \int dt p(\vec{x}, t) (E[t|\vec{x}] - t) +$$

2nd term vanishes

$$E[t|\vec{x}] p(\vec{x}) - p(\vec{x}) E[t|\vec{x}] = 0$$

$$+ \int dt d\vec{x} p(\vec{x}, t) (E[t|\vec{x}] - t)^2$$

Only the 1st term depends on  $y(\vec{x})$ , and will be = 0 (at min) iff

$$y(\vec{x}) = E[t|\vec{x}], \text{ consistent with (+)}$$



The 3rd term represents noise in the target values & is indep. of the model:

$$\int d\vec{x} p(\vec{x}) E[t|\vec{x}]^2 - 2 \int dt d\vec{x} t p(\vec{x}, t) E[t, \vec{x}] +$$

~~$$\int dt t^2 p(t)$$~~

$$+ \int dt t^2 p(t) =$$

$$= \int d\vec{x} \frac{\int dt t p(\vec{x}, t) \int dt' t' p(\vec{x}, t')}{p(\vec{x})} - \langle t^2 \rangle$$

$$p(t|\vec{x}) = \frac{p(\vec{x}, t)}{p(\vec{x})} \rightarrow -2 \int dt d\vec{x} t p(\vec{x}, t) \frac{\int dt' t' p(\vec{x}, t')}{p(\vec{x})} + \int dt t^2 p(t) =$$

$$= \langle t^2 \rangle - \int \frac{d\vec{x}}{p(\vec{x})} \int dt t p(\vec{x}, t) \int dt' t' p(\vec{x}, t')$$

$$\int d\vec{x} p(\vec{x}) \left( \frac{\int dt t p(\vec{x}, t)}{\int dt p(\vec{x}, t)} \right)^2 = \int d\vec{x} p(\vec{x}) \langle t(\vec{x}) \rangle^2$$

$$\langle t^2 \rangle = \int dt t^2 p(t) = \frac{\int dt d\vec{x} t^2 p(\vec{x}, t)}{\int dt d\vec{x} p(\vec{x}, t)} \stackrel{\text{"1"}}{\parallel} \int d\vec{x} \langle t^2(\vec{x}) \rangle p(\vec{x})$$

$\langle \dots \rangle$  signifies an average wrt  $p(\vec{x}, t)$

$$\langle t^2(\vec{x}) \rangle = \frac{\int dt t^2 p(\vec{x}, t)}{\int dt p(\vec{x}, t)} = \frac{\int dt t^2 p(\vec{x}, t)}{p(\vec{x})}$$

$$\text{So, } \int dt d\vec{x} p(\vec{x}, t) (E[t|\vec{x}] - t)^2 = \int d\vec{x} p(\vec{x}) [\langle t^2(\vec{x}) \rangle - \langle t(\vec{x}) \rangle^2]$$