

$$w_{k+1, j} = w_{k, j} - \alpha \left[\frac{\partial l(w_k)}{\partial w_{k, j}} \right]$$

$$\frac{\partial l(w_k)}{\partial w_{k, j}} = \frac{\partial}{\partial w_{k, j}} \frac{1}{2} \sum_{i=0}^n (y_i - \hat{y}_i)^2$$

$$= \sum_{i=0}^n \frac{\partial}{\partial w_{k, j}} \frac{1}{2} (y_i - \hat{y}_i)^2$$

$$= \sum_{i=0}^n \cancel{1} (y_i - \hat{y}_i) \frac{\partial}{\partial w_{k, j}} (y_i - \hat{y}_i)$$

$$= \cancel{\sum_{i=0}^n} (y_i - \hat{y}_i) \frac{\partial}{\partial w_{k, j}} \sum_{\beta=1}^3 w_{k, \beta} x_{i, \beta}$$

$$= \cancel{\sum_{i=0}^n} (y_i - \hat{y}_i) \frac{\partial}{\partial w_{k, j}} \underbrace{\sum_{\beta=1}^3 w_{k, \beta} x_{i, \beta}}_{z \cdot \text{const} = \text{const}}$$

$$= \cancel{\sum_{i=0}^n} (y_i - \hat{y}_i) x_{i, j}$$

skipping
steps

$$l(w_k) = \frac{1}{2} \sum_{i=0}^n (y_i - f_{w_k}(x_i))^2$$

$$f_w(x_i) = \sum_{j=1}^3 w_j x_{ij}$$

$$\frac{\partial}{\partial x} x^2 \rightarrow 2x$$

$$\frac{\partial}{\partial x} f(x)^2 \rightarrow 2f(x) \frac{\partial f(x)}{\partial x}$$

what if $x_{i, j}$ is magnitude
 $\sim 10^{10}$?

$$f(x) = x^2$$

$$\nabla f(x) = 2x$$

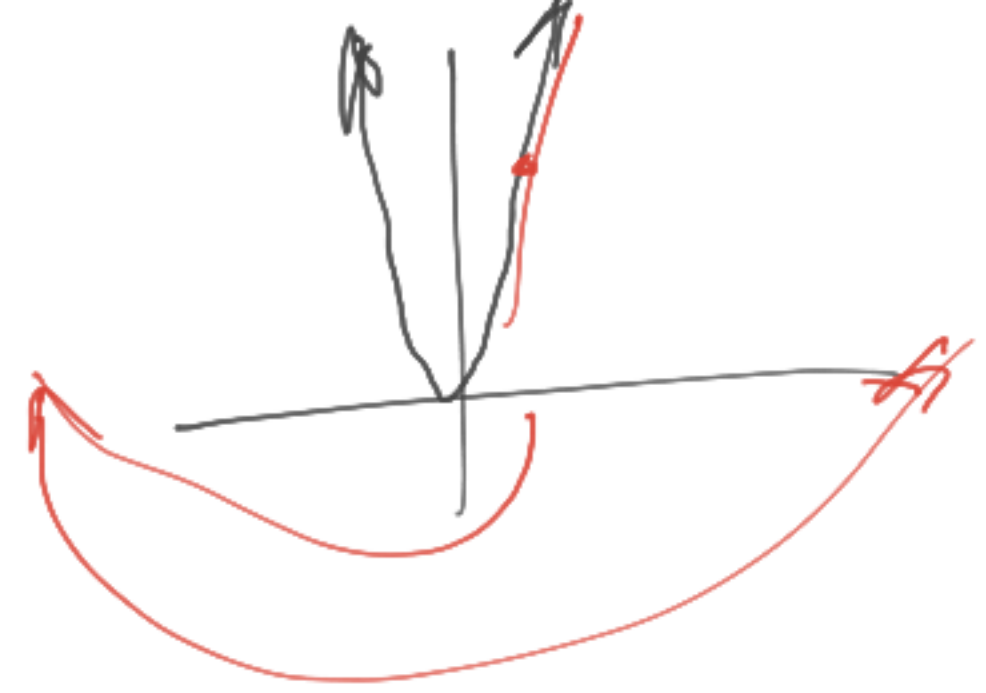
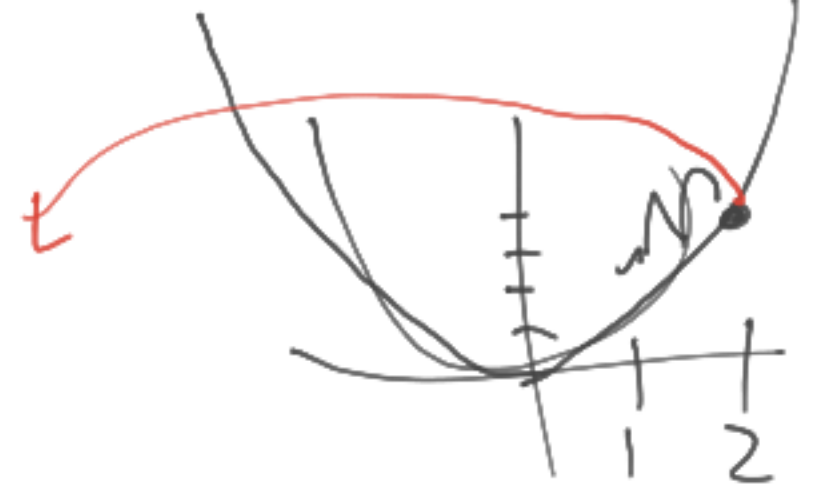
$$x_0 = 2$$

$$\alpha = .0001$$

$$x_1 = 2 - .0001 \cdot 2(2)$$

$$= 2 - .0004$$

$$= 1.9996$$



$\approx 2 - 2 \cdot 10^{-7} (4)$
massive negative number.

Input normalization

$$x'_{i,j} \in [-1, 1]$$

Before $x_{i,j} \in [a, b]$

$$x'_{i,j} = 2 \frac{x_{i,j} - a}{(b-a)} - 1$$

$$x'_{i,j} \in [0, 1]$$

$$x'_{i,j} = \frac{x_{i,j} - a}{(b-a)}$$

$$f(x) = 10^{10} x^2$$

$$\nabla f(x) = 2 \cdot 10^{10} x$$

$$x_0 = 2$$

$$\alpha = .0001$$

$$x_1 = 2 - .0001 \cdot 2 \cdot 10^{10} (4)$$

$$x_{i,1} \sim 10^{10}$$

$$\alpha \approx \frac{1}{10^{10}}$$

$$x_{i,2} \sim 1$$

step size is too small for $w_{i,2}$

too big step
 $\alpha = \frac{1}{10}$

$$\frac{\partial l(w_k)}{\partial w_{k,j}} = \sum_{i=1}^n (y_i - \hat{y}_i) x_{i,j}$$

Normalize:

- Mean zero

- Standard deviation 1
(variance 1)

$$X'_{i,j} = X_{i,j} - \frac{1}{n} \sum_{\beta=1}^n X_{\beta,j}$$

mean $X_{\cdot,j}$
zero.

		X		
		1	2	3
1 - X_1 :	7	14	31	$m=3$
2 - X_2 :	-2	-1	16	
3 - X_3 :	1	2	3	
4 - X_4 :	7	6	5	
5 - X_5 :	12	14	5	

$X_{5,2}$

$$\sigma^2 = \text{Var}(X_{\cdot,j})$$

$$X'_{i,j} = \frac{X_{i,j}}{\sqrt{\text{Var}(X_{\cdot,j})}}$$

$$\text{Var}(X'_{\cdot,j}) = \text{Var}\left(\frac{X_{\cdot,j}}{\sqrt{\text{Var}(X_{\cdot,j})}}\right)$$

$$= \frac{1}{\sigma^2} \text{Var}(X_{\cdot,j})$$

$$= \frac{1}{\sigma^2} \sigma^2 = 1$$

$$\text{Var}(aX) = a^2 \text{Var}(X) \quad \text{Var}\left(\frac{a}{c} X\right) = \frac{a^2}{c^2} \text{Var}(X)$$

$$\text{Var}(x_{1,j}, x_{2,j}, \dots, x_{n,j}) =$$

! sample

$$\frac{1}{n-1} \sum_{i=1}^n (x_{i,j} - \bar{x}_j)^2$$

$$\bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{i,j}$$

$H_{i,j}$

$$X_{i,j} = X_{i,j} - \frac{1}{n} \sum_{\beta=1}^n X_{\beta,j}$$

For each j

$$\bar{X}_j$$

$\sigma_j^2 = \text{variance}$

$$\bar{X}_j = \frac{1}{n} \sum_{i=1}^n X_{i,j}$$

$$\sigma_j^2 = \frac{1}{n-1} \sum_{i=1}^n (X_{i,j} - \bar{X}_j)^2$$

$H_{i,j} \rightarrow$

$$X_{i,j} = \frac{X_{i,j} - \bar{X}_j}{\sqrt{\sigma_j^2}}$$

For $j = 1$ to m

$$\bar{X}_j = \frac{1}{n} \sum_{i=1}^n X_{i,j}$$

$$\sigma_j^2 = \frac{1}{n-1} \sum_{i=1}^n (X_{i,j} - \bar{X}_j)^2$$

Bessel's correction

For $i = 1$ to n

For $j = 1$ to m

$$X_{i,j} \leftarrow \frac{X_{i,j} - \bar{X}_j}{\sqrt{\sigma_j^2}}$$

$$l(w_k) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$y_i \in [0, 4]$$

$$l(w_k) = 10 \mu F Z$$

~~X~~ $n=1$

$n=100,000$

1) Misscaled objectives hard to interpret.

2) Misscaled objectives hard to set α .

\sim scale (1)

\times scale (100,000)

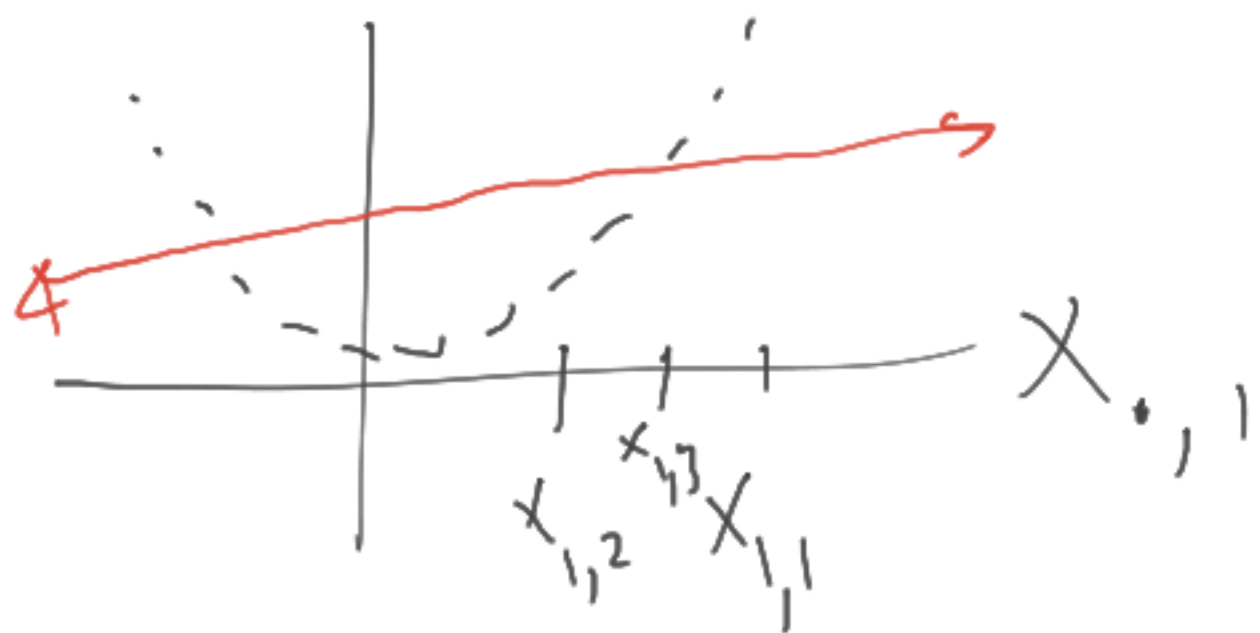
1) Check scaling of y_i .

2) Divide by n .

Basis functions.

$m=1$

"Affine"



Linear or Affine means

$\frac{\partial f_w(x_i)}{\partial w} = \text{constant}$.
Cost function of w .

"linear"

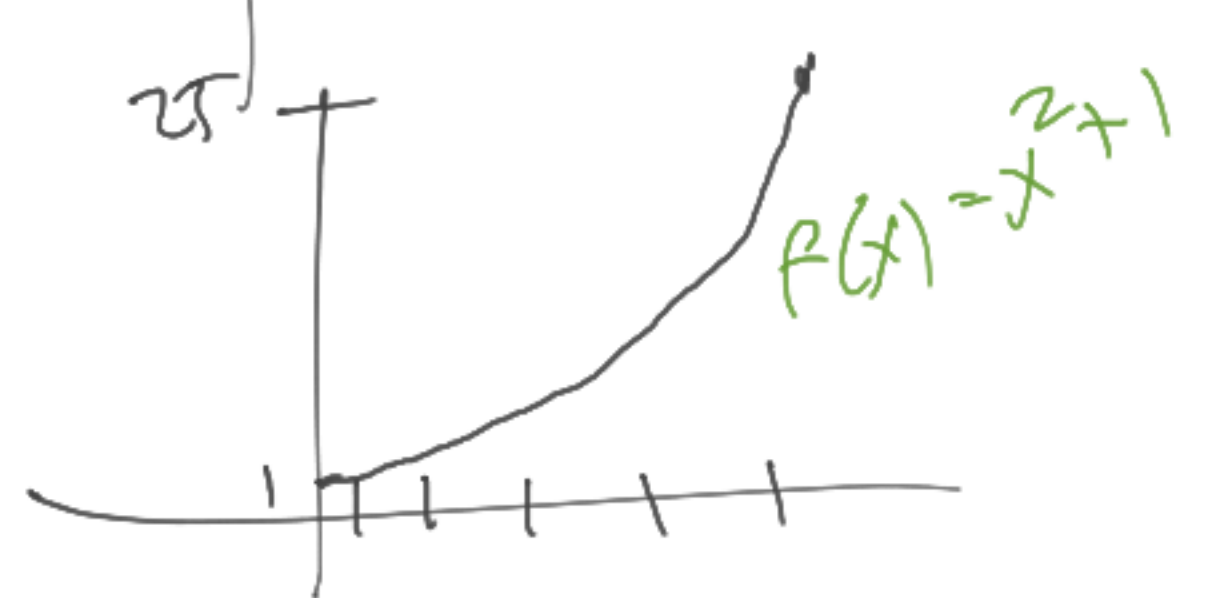
$$f_w(x_i) = \sum_{j=1}^m w_j x_{i,j}$$

↓
linear model
is linear w.r.t.
 w .

$$\frac{\partial}{\partial w_\beta} f_w(x_i) = \sum_{i=1}^n x_{i,\beta}$$

$X, Y \rightarrow X', Y$
 $m=1 \quad m'=2$
 $x^2, 1$

1	2	1	1	2
2	5	4	1	5
3	10	9	1	10
4	17	16	1	17
5	26	25	1	26



$w_1 = 1$
 $w_2 = 1$

$f_w(x_i) = w_1 x'_{i,1} + w_2 x'_{i,2}$
 $= w_1 x^2 + w_2 1$
 $= x^2 + 1$

ϕ "phi"

phi-features

features.

$$f_w(x_i) = \sum_{j=1}^m x_{i,j} w_j$$

$$x'_i = \phi(x_i)$$

$$\phi: \mathbb{R}^m \rightarrow \mathbb{R}^{m'}$$

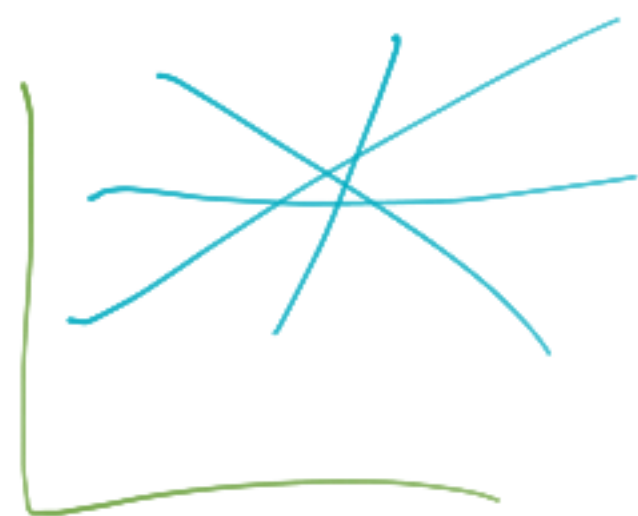
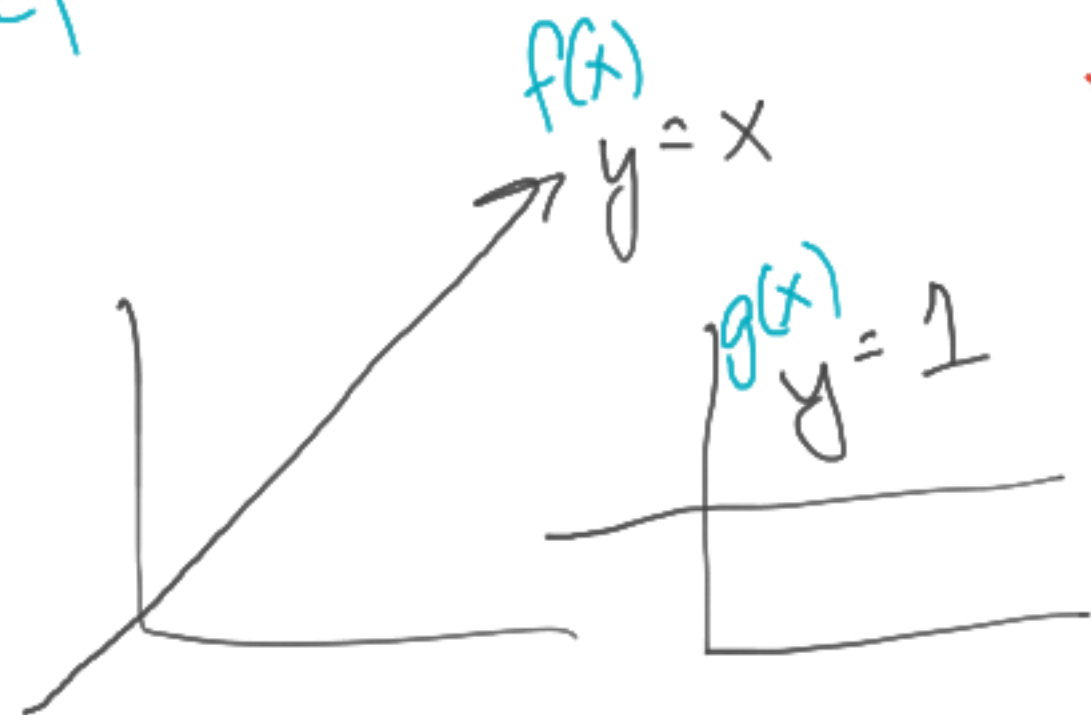
$$f_w(x_i) = \sum_{j=1}^m \phi(x_{i,j}) w_j$$

vector of features

Hyperparameter.

"basis function"

$$c_1 f(x) + c_2 g(x)$$



How use features from \mathbb{R}^2



$$+ \begin{pmatrix} c_1 [0, 1] \\ c_2 [1, 0] \end{pmatrix}$$

Basis for \mathbb{R}^2

Polynomial basis. (monomial basis)

$x_{i,1}$

$x_{i,2}$

$x_{i,3}$

$m=3$

m'

"order"

Hyperparameter.

\hookrightarrow

$x_{i,1}$

$x_{i,1}^2$

$x_{i,1}^3$

\dots

$x_{i,1}^0$

$x_{i,2}$

$x_{i,2}^2$

\dots

$x_{i,2}^0$

\dots

(No terms with both $x_{i,j}$ and $x_{i,j}'$)

$x_{i,1}^4$ $x_{i,2}^3$

\rightarrow Fourier Basis.

\rightarrow cos sin.

- Normalize

before

applying
basis.

