

④ Supervised learning I

* Supervised: $\mathcal{D} = \{\tilde{x}^i, \tilde{y}^i\}$ * Example $\left\{ \begin{array}{l} y = \text{temperature} \\ x_1 = \text{day}; \\ x_2 = \text{place}; \end{array} \right.$
 features outputs

a) Linear regression and Dynamic regression

* $\min_{\theta} \ell(\theta) = \text{MSE}(\theta) = \|y - X\theta\|_2^2$; The output $y \in \mathbb{R}^N$
 where $x = [1, x_1, x_2, \dots]$ for one feature and $X = [1, x_1, x_2, \dots, 1, x_1, x_2, \dots]$

* We can use different norms / or penalty functions

* Regularization

* L^2 regularization (Tikhonov regularization / Ridge regression)

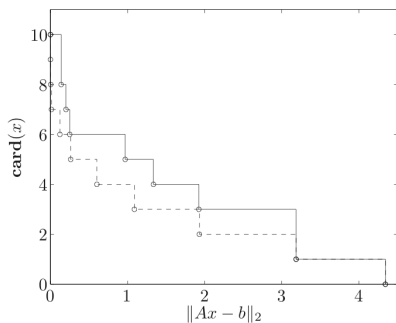
$$\min \|y - X\theta\|_2^2 + \lambda \underbrace{\|\theta\|_2^2}_{\theta^T \theta}$$

$$\underbrace{(X^T X + \lambda I)}_{\text{full rank}} \theta = X^T y$$

- $\lambda > 0$
- Prior knowledge of θ
- Penalize large values of θ
- Also useful when X is square but bad conditioned i.e. Elasticity?

* Other possible regularization

* L^1 regularization $|\theta|$ (Lasso problem)



$$\min \|y - X\theta\|_2^2$$

$$\text{s.t. } \text{card}(\theta) \leq K$$

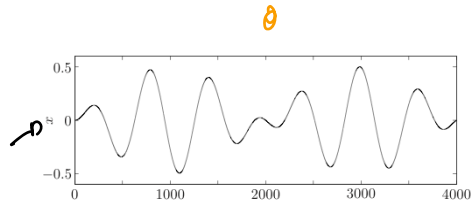
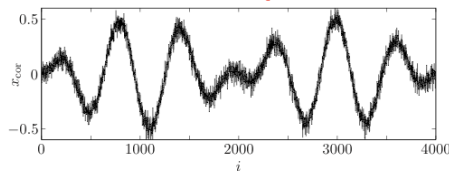
asking for a sparse representation of θ

good approximation

$$\min \|y - X\theta\|_2^2 + \lambda \|\theta\|_1$$

* $\|\nabla \theta\|^2 \rightarrow$ smoothing (Laplacian)

$$\min \|y - I\theta\|_2^2 + \lambda \|\nabla \theta\|^2$$



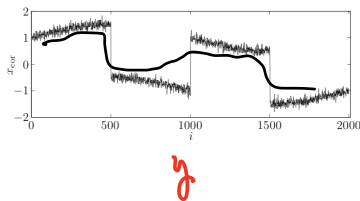
In variational form: $\min \int (y - \theta)^2 + \epsilon \int \nabla \theta^2 = f(\theta)$

Euler-Lagrange derivative: $\frac{\delta f(\theta)}{\delta \theta}(r) = 2 \int (y - \theta) \psi + \epsilon \int \nabla \theta \nabla \psi = 0 \quad \forall \psi$

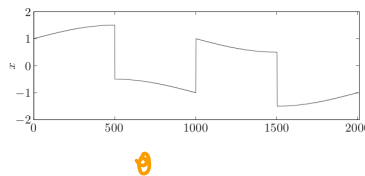
$$\left(\int \psi + \epsilon \int \nabla \psi \nabla \theta = \int y \psi \quad \forall \psi \Rightarrow \begin{cases} \epsilon \Delta \theta + \theta = y \\ [1 + \epsilon k] \theta = \bar{y} \end{cases} \right. \quad \text{Strong form}$$

* $|\nabla \theta| \rightarrow$ Total variation

$$\min \|y - I\theta\|_2^2 + \lambda |\nabla \theta|$$



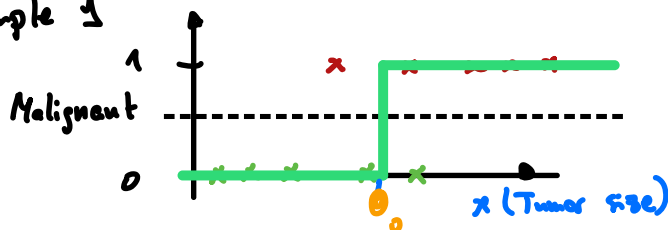
θ^*
Denoising



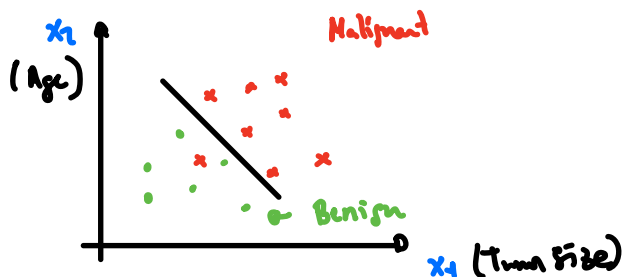
b) Logistic regression / classification

* Classification: the output variable $y \in \{0, 1\}$

* Example 1



* Example 2



* Hypothesis space / hypothesis function

$$h_{\theta}(x) = x\theta$$

* Thresholding (Heaviside function)

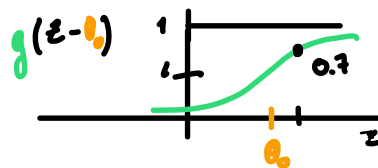
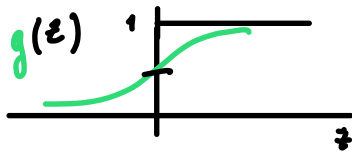
$$g(h_{\theta}(x)) = \begin{cases} 0 & h_{\theta}(x) < 0 \\ 1 & h_{\theta}(x) \geq 0 \end{cases}$$

thresholding



* Sigmoid function: (Logistic function)

$$g(z) = \frac{1}{1 + e^{-z}}$$



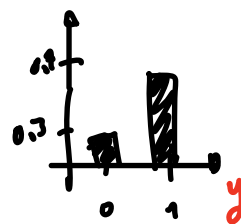
* Interpretation of $g(h_{\theta}(x))$:

Estimated probability that $y=1$ on input x

For some x_1 so $g(h_{\theta}(x_1)) > 0.7 \Rightarrow 70\%$ chance of tumor being malignant

$$g(h_{\theta}(x)) = p(y=1 | X=x; \theta)$$

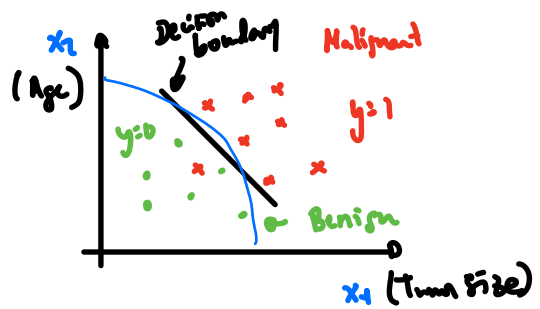
$$1 - g(h_{\theta}(x)) = 1 - p(y=1 | X=x; \theta) = p(y=0 | X=x; \theta) = 0.3$$



* Decision boundary

$$g(\theta^T x) = \begin{cases} 1 & \theta^T x \geq 0 \\ 0 & \theta^T x < 0 \end{cases}$$

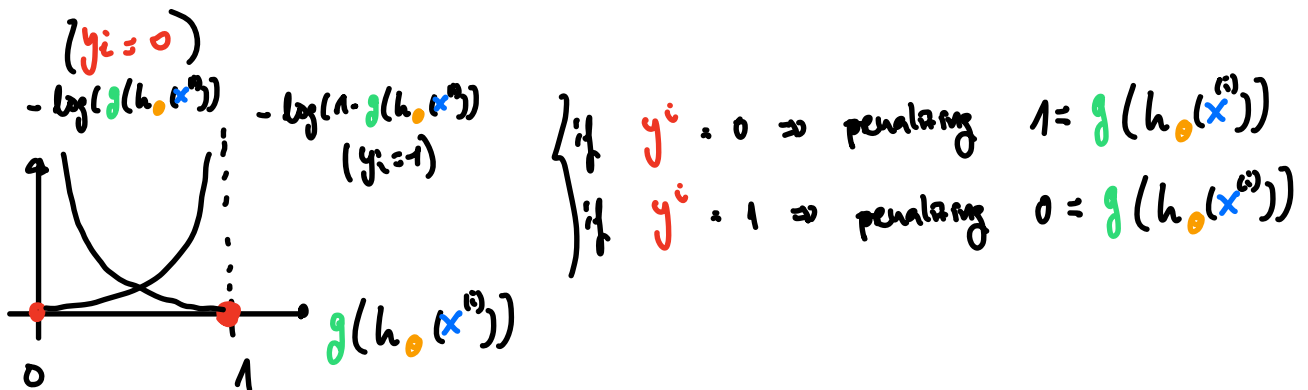
f_{linear}



* Cost function: Negative loglikelihood / cross-entropy

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left((1 - y^i) (-\log(1 - g(h_\theta(x^i)))) + y^i (-\log(g(h_\theta(x^i)))) \right)$$

$y^i = 0 \Rightarrow g = 1 \Rightarrow -\log(1 - g) \rightarrow \infty$



* Property of sigmoid function

$$g(z) = \frac{1}{1 + e^{-z}} ; g'(z) = g(z)(1 - g(z))$$

* Gradient

$$\frac{\partial J_0(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (1 - y^i) \frac{-1}{(1 - g(h_{\theta}(x^i)))} - (g(h_{\theta}(x^i)) (1 - g(h_{\theta}(x^i))) x^{(i)}_j) \\ + (y^i) \frac{-1}{(g(h_{\theta}(x^i)))} - (g(h_{\theta}(x^i)) (1 - g(h_{\theta}(x^i))) x^{(i)}_j)$$

$$= - \frac{1}{m} \sum_{i=1}^m [(1 - y^i) (g(h_{\theta}(x^i))) - (y^i) (1 - g(h_{\theta}(x^i)))] x^{(i)}_j$$

$$= - \frac{1}{m} \sum_{i=1}^m [g(h_{\theta}(x^i)) - y^i] x^{(i)}_j = 0$$

the gradient is similar to polynomial regression

$$\min_{\theta} \frac{1}{2} \|x_0 - y\|^2 = f_{\theta}(\theta)$$

$$\frac{\partial f_{\theta}(\theta)}{\partial \theta} = x^T x \theta - x^T y = 0 \quad \nabla f_{\theta}(\theta) = 0 \\ = x^T [x_0 - y]$$

* Remarks

* You can use gradient method: $\theta_{k+1} = \theta_k - \alpha \nabla J_0(\theta_k)$

* Stochastic gradient, randomly select L values of the data $L \leq m$

$$\nabla J_0(\theta_k) \approx \frac{1}{L} \sum_{i=1}^L [g(h_{\theta}(x^i)) - y^i] x^{(i)}$$

* Regularization also important

Typically we take hundreds of data to compute the stochastic gradient. When we have used all the data we say that we spend an epoch.

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^m (1 - y^i) (-\log(1 - g(h_{\theta}(x^i)))) + y^i (-\log(g(h_{\theta}(x^i)))) \\ + \lambda \theta^T \theta$$

For L_1 norm regularization use: Proximal methods

$$\frac{\partial J_0(\theta)}{\partial \theta} = \frac{1}{m} \sum_{i=1}^m [g(h_{\theta}(x^i)) - y^i] x^{(i)} + \lambda \theta$$

- * Gradient expensive \rightarrow Stochastic
- * Stopping criteria is expensive \rightarrow $\| \nabla J \|, \dots$ n iterations
- * Computing cost expensive \rightarrow We don't compute it
 - \hookrightarrow Exact line-search
 - \hookrightarrow Backtracking \rightarrow α very small but fixed

