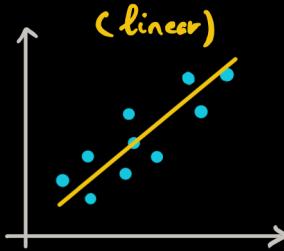
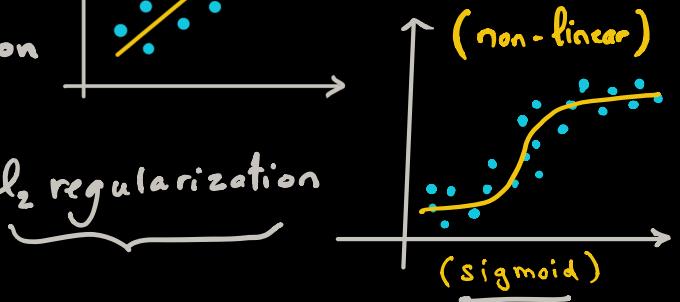


MML6

1. Basics of linear regression



2. Non-linear regression and ℓ_2 regularization

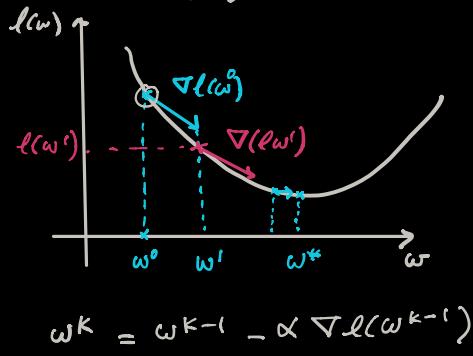


3. Feature design for linear regression

Max \mathcal{L}_5 ?

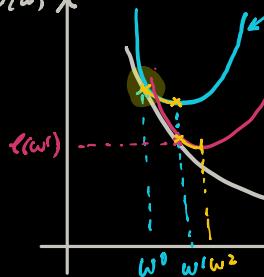
Gradient descent

- α : learning rate
- w^0 : initial guess
- ε : stopping "value"



Newton's method

- w^0
- ε ; $\|\nabla l(w^*)\| \leq \varepsilon$



$$w^K = w^{K-1} - \underbrace{\left[\nabla^2 l(w^{K-1}) \right]^{-1}}_{\text{computationally expensive}} \nabla l(w^{K-1})$$

- ⊕ if l is convex \Rightarrow faster convergence
- ⊕ no α
- ⊖ if l is non-convex \Rightarrow not converge
- ⊖ $\left[\nabla^2 l(w) \right]^{-1}$ computationally expensive

Estimating a predictive learner { classifier
regressor }

Think about ↵ ① define loss function & its set of parameters
how to define $\underline{l(w)}$

optimization ↵ ② minimize $\underline{l(w)}$ using the training data

cross-validation ↵ ③ Once the optimal set of model parameters w^* is estimated, evaluate model on testing data.

Basics of linear regression

Training step

Training Samples:

$$\{(x^1, y^1), (x^2, y^2), \dots, (x^n, y^n)\} = \{\tilde{x}_i, \tilde{y}_i\}_{i=1}^n$$

\tilde{x}_i feature vector (input) $\in \mathbb{R}^d$

\tilde{y}_i scalar-valued target (output) $\in \mathbb{R}$

Case $d=1$ (scalar input)

$$y = \tilde{x}^P \tilde{w} + b$$

Slope "y-intercept"

* what is the approximation or prediction of the target y^P by this model?

Case $d \geq 1$ (general case)

$$\tilde{x}^P = \begin{bmatrix} \tilde{x}_1^P \\ \tilde{x}_2^P \\ \vdots \\ \tilde{x}_d^P \end{bmatrix} \Rightarrow \left\{ \begin{array}{l} \tilde{y}^P \approx \tilde{x}^P \tilde{w} + b \\ p = 1, \dots, n \end{array} \right.$$

weight

bias

$$\tilde{y}^P \approx \tilde{x}^P \tilde{w} + b$$

$\tilde{w} \in \mathbb{R}^{d+1}$

$f(\tilde{w}, b) = f_{\tilde{w}, b}(\tilde{x}) \approx \tilde{y}^P$

\tilde{x}^P point $(\tilde{x}^P, \tilde{y}^P)$ lies close to the hyperplane $y = \tilde{x}^P \tilde{w} + b$.

Example

Guggenmos et al. *Translational Psychiatry* (2017) 2:1279
DOI 10.1038/s41398-017-0037-y

Translational Psychiatry

ARTICLE

Open Access

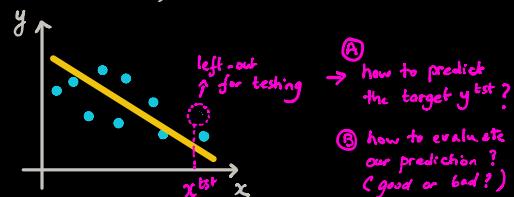
Quantitative neurobiological evidence for accelerated brain aging in alcohol dependence

Matthias Guggenmos¹, Katharina Schmack¹, Maria Sekutowicz¹, Maria Garbusow¹, Miriam Sebold¹, Christian Sommer², Michael N. Smolka^{2,3}, Hans-Ulrich Wittchen⁴, Ulrich S. Zimmermann², Andreas Heinz¹ and Philipp Sterzer¹

Testing Step

Given an available dataset $\{(x_i, y_i)\}_{i=1}^n$

use leave-one-out cross validation to train our linear regressor:



$$\tilde{w}^* = \left[\sum_{i=1}^{n-1} (\tilde{x}_i \tilde{x}_i^T) \right]^{-1} \sum_{i=1}^{n-1} \tilde{x}_i \tilde{y}_i$$

$\tilde{x}_i^* = \begin{bmatrix} 1 \\ \tilde{x}_i \end{bmatrix}$; $\tilde{w} = \begin{bmatrix} b \\ w \end{bmatrix}$

⑥ how to predict the target y^{tst} ? [step-by-step]

$$\hat{y}^{\text{tst}} = \tilde{x}^{\text{tst}} \tilde{w}^* + b^*$$

(n-1) samples

⑦ how to evaluate our prediction? (good or bad?)

Given the optimal solution:

$$\tilde{w}^* = \begin{bmatrix} b^* \\ \tilde{w}^* \end{bmatrix}$$

evaluate our prediction for each left-out subject using

$$\text{TSE} = \text{Total Squared Error}$$

$$\text{TSE} = \frac{1}{n} \sum_{i=1}^n (\tilde{x}_i^T \tilde{w}^* - \tilde{y}_i)^2$$

$$\tilde{w}^* = \left[\sum_{i=1}^n \tilde{x}_i \tilde{x}_i^T \right]^{-1} \left[\sum_{i=1}^n \tilde{x}_i \tilde{y}_i \right]$$

what if we leave 3 subjects out? Would this change the value we get for TSE?

$$\text{MSE}_{100} \neq \text{MSE}_{5\text{-fold}}$$

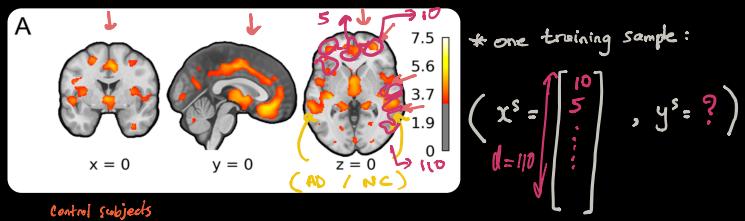
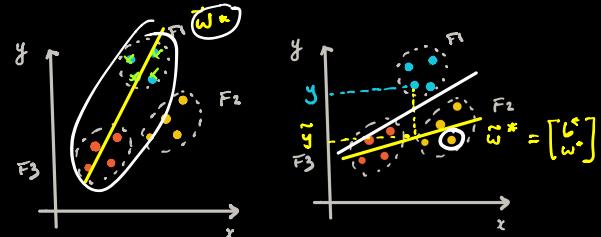
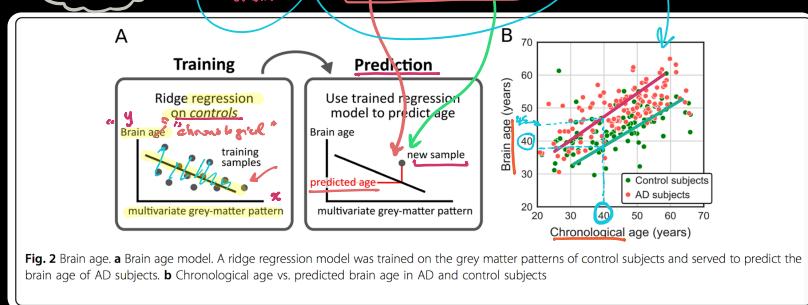


Fig. 1 Correspondence between AD-related and age-related grey matter loss (GML). **a** and **b** show t-maps for univariate whole-brain analyses, thresholded at $p < 0.001$ uncorrected, for illustration. **a** T-map for AD-related grey matter volume loss based on a two-sample t test between AD subjects and control subjects, controlling for age, gender, site, smoking (FTND sum score) and general health status (WHODAS).



K-fold CV:

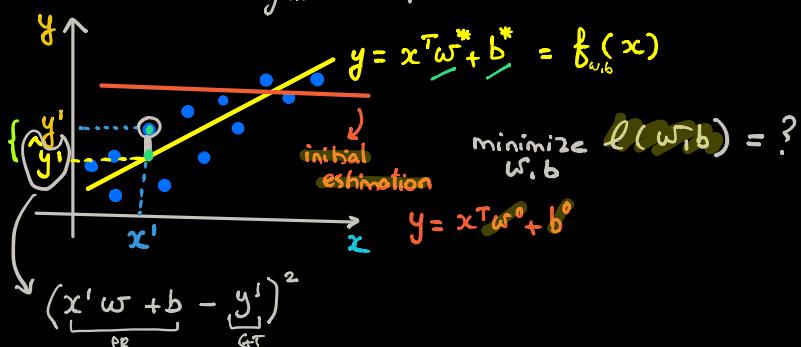
$\text{overall MSE} = \frac{\text{MSE}(F1) + \text{MSE}(F2) + \dots + \text{MSE}(FK)}{K}$

② Least squares loss function for linear regression

$$\arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{(\mathbf{x}^i, y^i) \in D} \underbrace{E(f_{\mathbf{w}}(\mathbf{x}^i), y^i)}_{\ell(w, b)} + R(\dots)$$

- Our goal \Rightarrow find the parameters of the hyper plane which best fits a regression dataset.
- Sol \Rightarrow these parameters $(w, b) \in \mathbb{R}^{d+1}$

Think of a simple loss function that we can minimize to find the optimal (w^*, b^*) .



$$\underset{b, w}{\text{minimize}} \sum_{i=1}^n \left(\underbrace{x_i^T w + b - y_i}_{\text{least squares}} \right)^2$$

$\underset{b, w}{\text{minimize}} \sum_{i=1}^n \frac{(x_i^T w + b - y_i)^2}{f(x_i^T w + b)}$
 sum of least squares

. Minimization of the loss function $\ell(w, b)$:

$$\tilde{x}_i^T \tilde{w} = ?$$

$$\tilde{x}_i \in \mathbb{R}^{d+1}, \quad \tilde{w} = \begin{bmatrix} b \\ w \end{bmatrix}$$

$$[1 \ x_i] \begin{bmatrix} b \\ w \end{bmatrix} = b + x_i^T w$$

→ using this compact notation where we have only one variable \tilde{w}
to minimize:

$$\ell(\tilde{w}) = \sum_{i=1}^n ((\tilde{x}_i^T \tilde{w} - y_i)^2)$$

→ compute the gradient of $\ell(\tilde{w})$:

$$\frac{df^2}{dx} = 2f'$$

$$\frac{\partial x^T w}{\partial w} = x$$

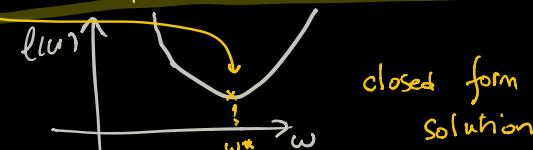
$$\nabla \ell(\tilde{w}) = \frac{\partial \ell(\tilde{w})}{\partial \tilde{w}} = \frac{\partial}{\partial \tilde{w}} \left[\sum_{i=1}^n ((\tilde{x}_i^T \tilde{w})^2 - 2\tilde{x}_i^T \tilde{w} y_i + y_i^2) \right]$$

$$= \sum_{i=1}^n [2\tilde{x}_i \tilde{x}_i^T \tilde{w} - 2\tilde{x}_i^T y_i]$$

$$= 0$$

$$\sum_{i=1}^n \tilde{x}_i \tilde{x}_i^T \tilde{w} = \sum_{i=1}^n \tilde{x}_i y_i$$

$$\tilde{w}^* = \left[\sum_{i=1}^n \tilde{x}_i \tilde{x}_i^T \right]^{-1} \sum_{i=1}^n \tilde{x}_i y_i$$



{ Algebra trick }

$$\left(\sum_{i=1}^n \tilde{x}_i \tilde{x}_i^T \right) \tilde{w} = \sum_{i=1}^n \tilde{x}_i y_i \quad \leftarrow \text{linear system to solve}$$

$$\Rightarrow \tilde{X} \tilde{X}^T \tilde{w} = \tilde{X} y$$

$\tilde{X} \tilde{y}$ (stack column wise) = $\sum_{i=1}^n \tilde{x}_i \tilde{y}^i$

The diagram illustrates the matrix-vector multiplication Xy . It shows a matrix X with columns $\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^n$ and a vector y with components y^1, y^2, \dots, y^n . The result is a scalar value. A green bracket underlines the term $\tilde{x}_i \tilde{y}^i$ in the expanded form $\tilde{x}_1 \tilde{y}^1 + \tilde{x}_2 \tilde{y}^2 + \dots + \tilde{x}_n \tilde{y}^n$, which is highlighted in green. The entire expression is enclosed in a green box.

$$\tilde{X} \tilde{X}^T \tilde{\omega} = \tilde{X} \tilde{y}$$

Algebraic solution to
the linear least squares loss
function

$$\ell(\tilde{\omega}) = \sum_{i=1}^n (\tilde{x}_i^T \tilde{\omega} - y^i)^2$$

$$\tilde{\omega}^* = [\tilde{X} \tilde{X}^T]^{-1} \tilde{X} \tilde{y}$$

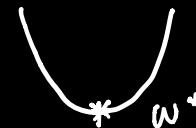
closed form sol^o if $\tilde{X} \tilde{X}^T$ is
invertible.

\therefore Still, better to solve the original linear
system as computing the inverse of
a matrix might be computationally
expensive.

②

Non-linear regression
an l_2 regularization

A) Previously, the loss function is convex & linear in w

$$\left\{ \begin{array}{l} \ell(w, b) = (x^T w + b - y)^2 \\ \text{or} \\ \ell(\tilde{w}) = (\tilde{x}^T \tilde{w} - y)^2 \quad \tilde{w} = \begin{bmatrix} b \\ w \end{bmatrix} \end{array} \right.$$


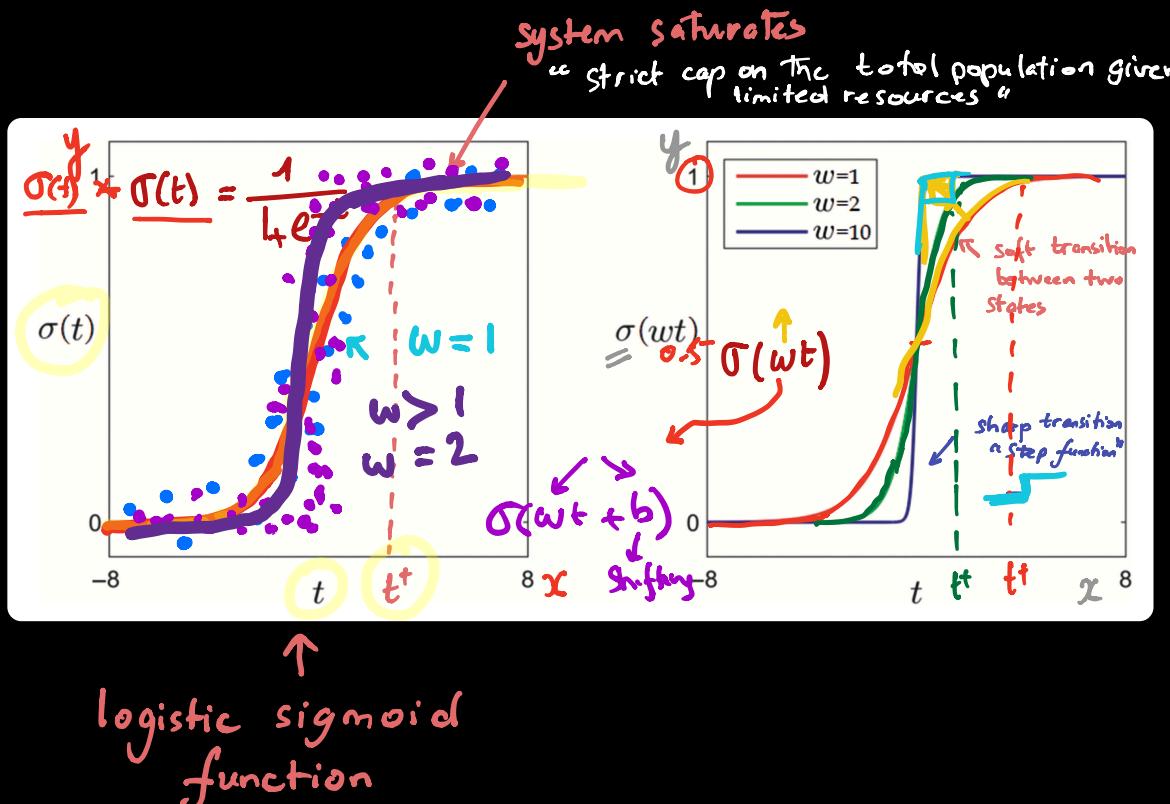
B) loss function is non-convex & nonlinear in w :



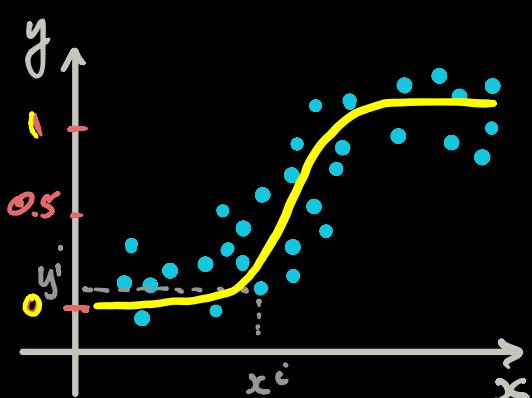
this requires an additional regularization term to improve loss optimization.

Example : Logistic regression

classic non-linear & non convex problem



- σ was designed to model population growth over time.
- At some point $t \geq t^+$, the growth stabilizes and levels off due to lack of limited resources (e.g., food, space, etc.).
- Given a dataset $\{(x_i, y_i)\}_{i=1}^n$ that is distributed like a sigmoid function :

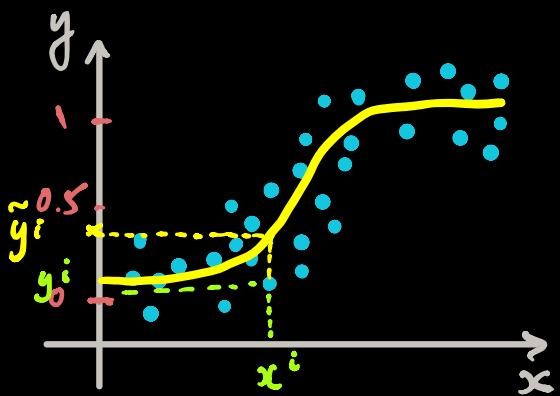


Goal = estimate the best sigmoid function fitting our data.

Sol^o = parameterize $\sigma(w, b)$

$x \rightarrow w x + b$

Smoothness of the transition $y=0 \rightarrow y=1$



The data satisfies

$$\tilde{y}^i = \sigma(x^i \omega + b) \approx y^i$$

\Downarrow

$$x \in \mathbb{R}^d$$

$$\sigma(\underbrace{x^T \omega}_{} + b) \approx y^i$$

\Downarrow

$$\underbrace{\sigma}_{\in \mathbb{R}^d}(\underbrace{x^T \omega}_{} + b) \approx y^i$$

↑
value
(shifting factor)

$$\ell(w, b) = \sum_{i=1}^n (\sigma(\underbrace{x^{i T} \omega}_{} + b) - \underbrace{y^i}_{})^2$$

training set: $\{(x^i, y^i)\}_{i=1}^n$

- Compact form:

$$\tilde{x}^i = \begin{bmatrix} 1 \\ x^i \end{bmatrix}, \quad \tilde{w} = \begin{bmatrix} b \\ \omega \end{bmatrix}$$

$$\ell(\tilde{w}) = \sum_{i=1}^n (\sigma(\underbrace{\tilde{x}^{i T} \tilde{w}}_f) - \underbrace{y^i}_{})^2$$

$$\frac{\partial x^T w}{\partial w} = x$$

$$\rightarrow \nabla \ell(\tilde{w}) = ? f$$

Hint: $\frac{\partial \sigma(x^T \omega)}{\partial \omega} = \underbrace{\sigma(x^T \omega)}_{\sigma'(x^T \omega)} (1 - \sigma(x^T \omega)) x$

$$\ell(\tilde{w}) = \sum_{i=1}^n \left[\underbrace{\sigma(\tilde{x}_i^\top \tilde{w})^2}_{\text{green}} - 2 \underbrace{\sigma(\tilde{x}_i^\top \tilde{w}) y_i}_{\text{green}} + y_i^2 \right]^0$$

$$\frac{\partial \ell(\tilde{w})}{\partial \tilde{w}} = \sum_{i=1}^n \left[\underbrace{2 \sigma'(\tilde{x}_i^\top \tilde{w}) (\sigma(\tilde{x}_i^\top \tilde{w}))}_{\text{green}} (1 - \sigma(\tilde{x}_i^\top \tilde{w})) \tilde{x}_i^i - 2 y_i \underbrace{\sigma'(\tilde{x}_i^\top \tilde{w})}_{\text{green}} (1 - \sigma(\tilde{x}_i^\top \tilde{w})) \tilde{x}_i^i \right]$$

$$= 2 \sum_{i=1}^n \left[(\sigma(\tilde{x}_i^\top \tilde{w}) - y_i) \underbrace{\sigma(\tilde{x}_i^\top \tilde{w}) (1 - \sigma(\tilde{x}_i^\top \tilde{w}))}_{\text{green}} \tilde{x}_i^i \right]$$

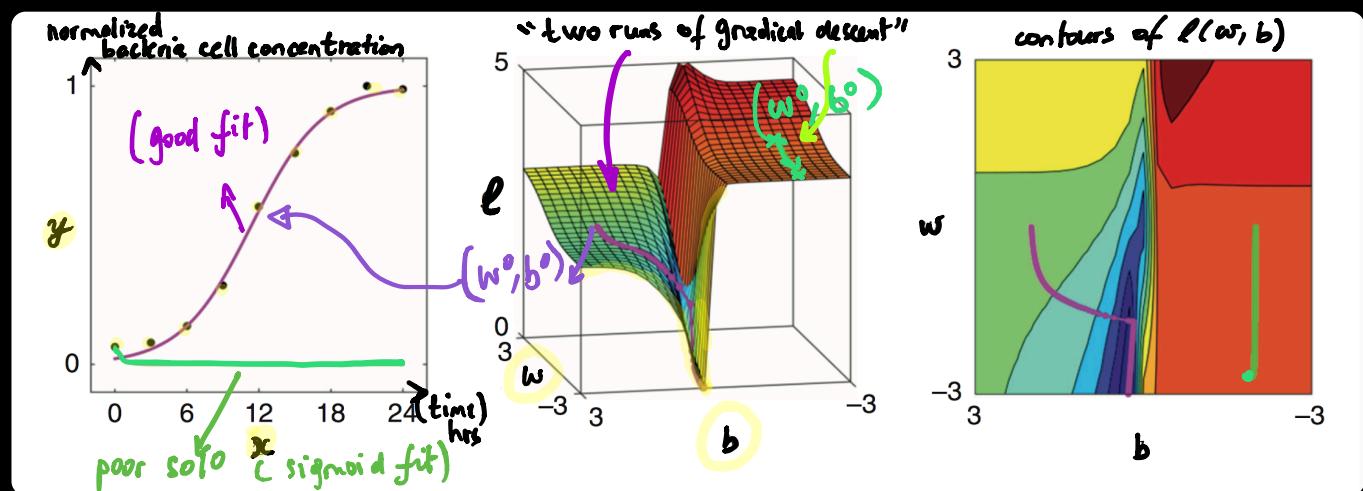
$$\Rightarrow \nabla \ell(\tilde{w}) = 2 \sum_{i=1}^n (\sigma(\tilde{x}_i^\top \tilde{w}) - y_i) \underbrace{\sigma(\tilde{x}_i^\top \tilde{w})}_{\text{green}} \underbrace{(1 - \sigma(\tilde{x}_i^\top \tilde{w}))}_{\text{green}} \tilde{x}_i^i$$

$\tilde{w}^* = []$

↑

Gradient descent can be used to optimize ℓ !

Example 1

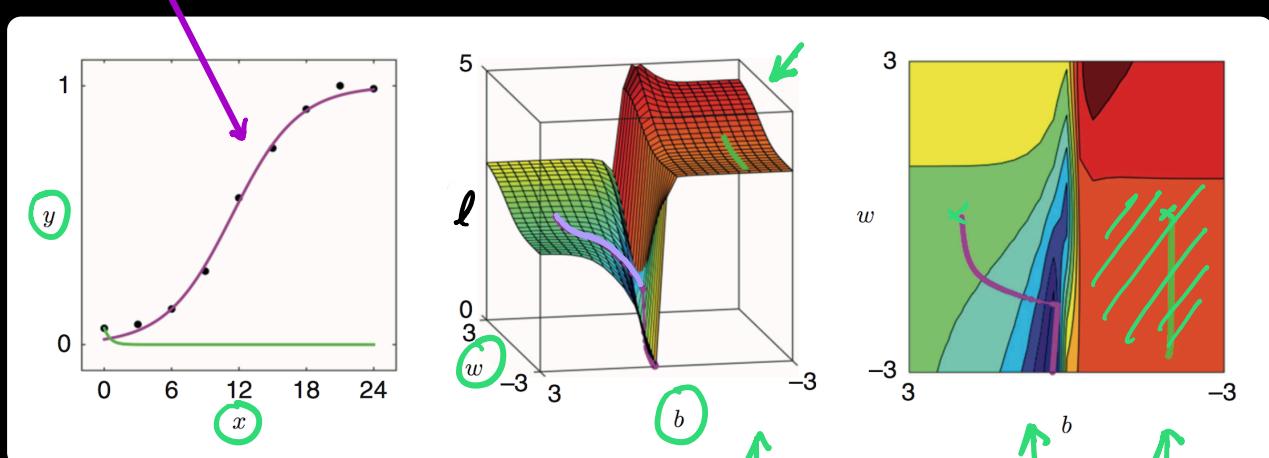


∴ our soln (w, b) depends on the initialization of the gradient descent algorithm.

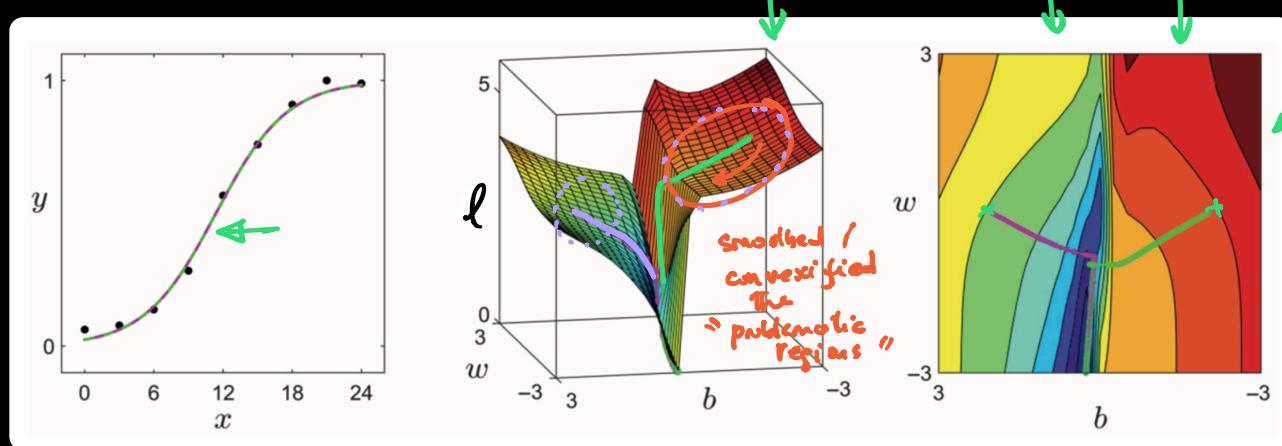
∴ initialization in the large flat orange region leads to a poor solution.

logistic

THINKING CAPS ON!

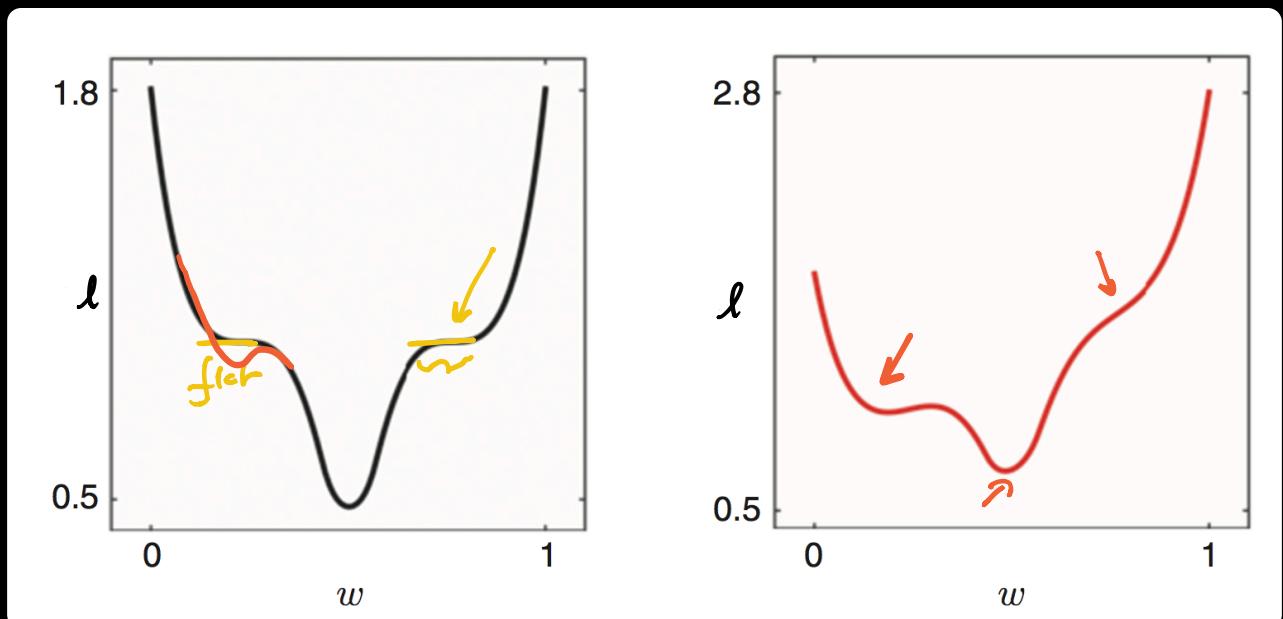


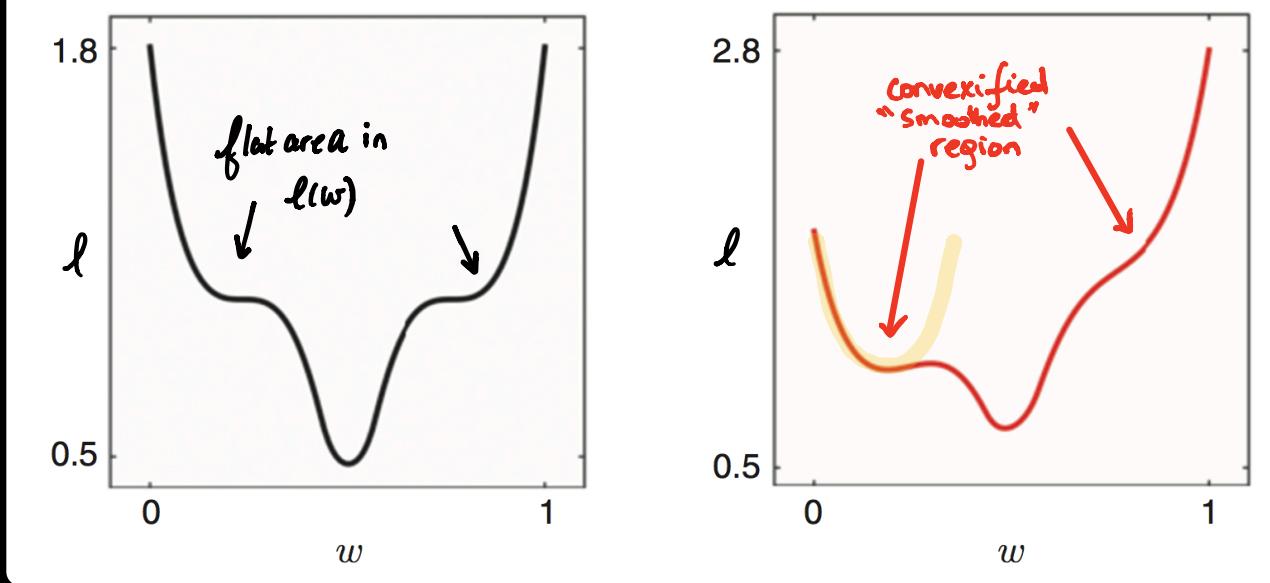
①



②

2D example





- $\ell(w, b)$ fit to data samples + $\lambda \|\mathbf{w}\|_2^2$ global ℓ_2 regularizer of the loss function
- $\|\mathbf{w}\|_2^2 = \sum_{i=1}^d w_i^2$

- if $\lambda = 0$ \Rightarrow only loss \Rightarrow trade-off parameter
- if λ very large \Rightarrow the regularizer draws out the loss function and we have

$$\ell(w, b) + \lambda \|\mathbf{w}\|_2^2 \approx \lambda \|\mathbf{w}\|_2^2$$

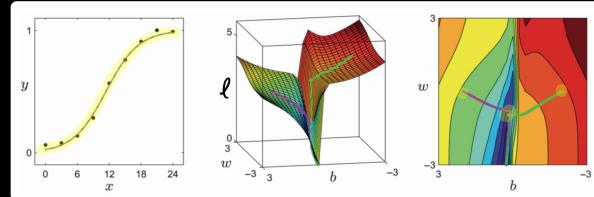
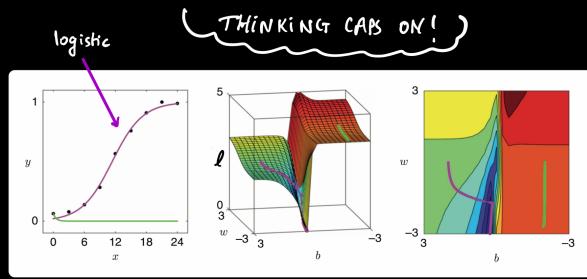
too large "takes over"
- Typically λ is set fairly small ($\lambda = 0.1$ or smaller)

Sol^o \Rightarrow

Improving problematic flat areas using ℓ_2 regularizer

$\ell_2 \|\cdot\|_2$ \Rightarrow a simple convex function that is added to the non-convex function to smooth it and convexify it in problematic non-convex regions.

* Logistic regression *



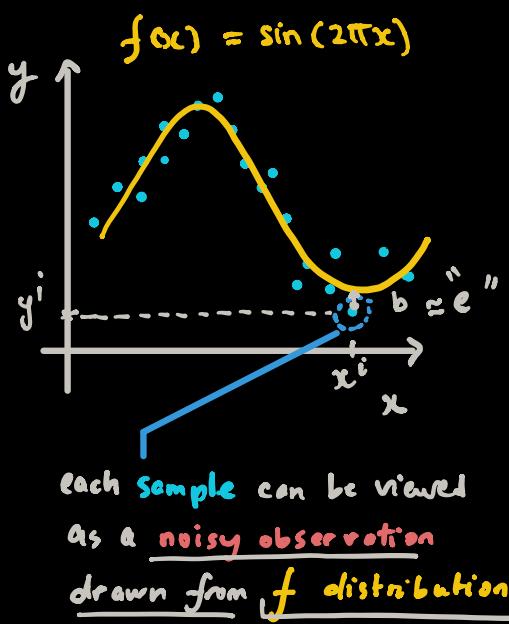
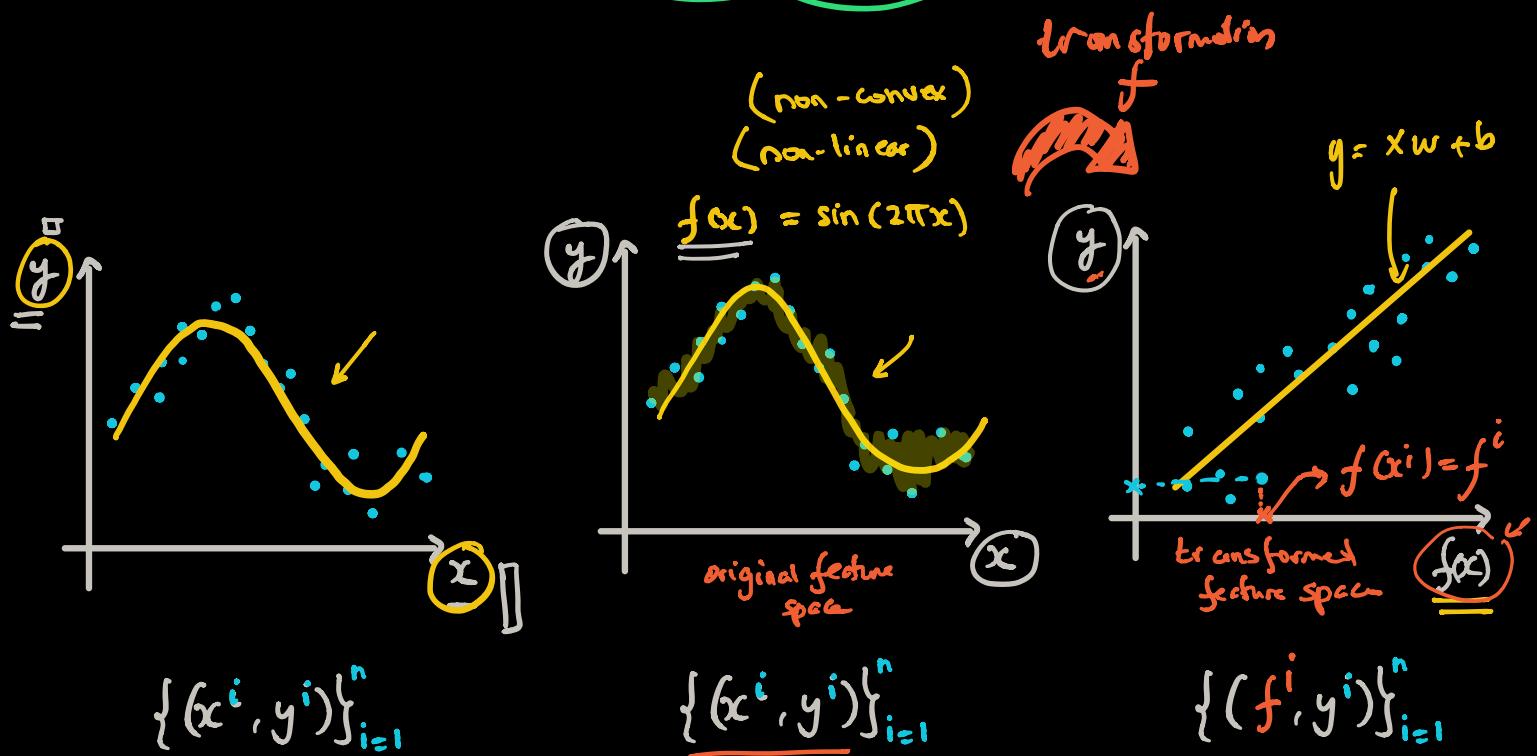
$$\ell(\omega, b) = \sum_{i=1}^n (\sigma(x_i^T \omega + b) - y^i)^2$$

$$\ell(\omega, b) = \sum_{i=1}^n (\sigma(x_i^T \omega + b) - y^i)^2 + \lambda \|\omega\|_2^2$$

$$\arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{(\mathbf{x}^i, y^i) \in \mathcal{D}} E(f_{\mathbf{w}}(\mathbf{x}^i, y^i)) + R(\dots)$$

Practice & code it up! exercise 3.11 exercise 3.13

③ Nonlinear relationship between input features and output values
 ↓
 Knowledge -driven feature design for regression



Knowledge-driven feature design
 our prior knowledge about the data
 ↓
 This might originate from our prior understanding / assumption about the data distribution.

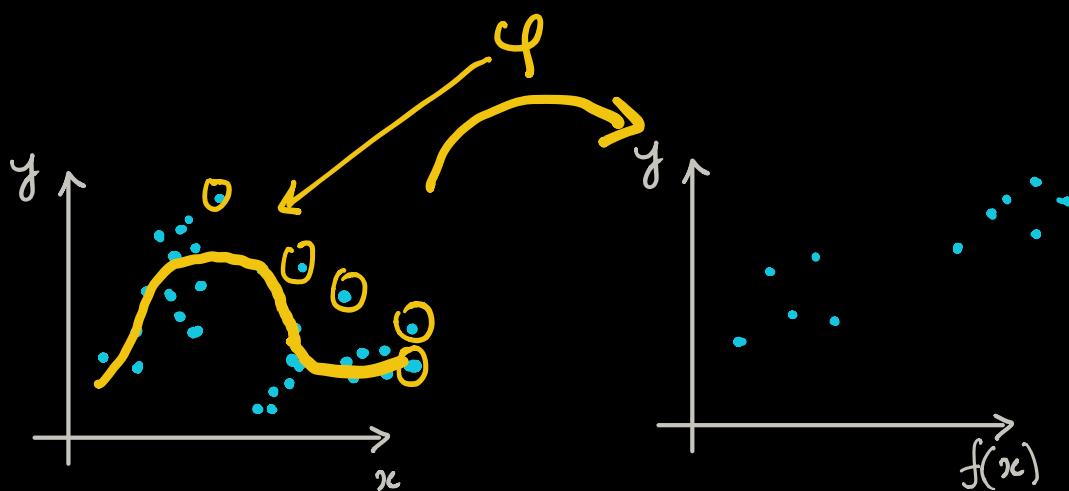
$$\frac{x^T}{\epsilon \|w\|} \rightarrow \frac{f(x^T)}{\epsilon \|w\|}$$

minimize b, w

$$\sum_{i=1}^n (f^i w - y^i)^2$$

\downarrow "The trick based on our knowledge of feature distribution in the original space"

mapped / transform feature



f is also called "estimated data generating function"
 (it is our estimated best fit to the dataset)

- Solve for w, b :

$$\underset{b, w}{\text{minimize}} \quad \sum_{i=1}^n (f_i^T w + b - y^i)^2$$

- Compact notation:

$$\tilde{f}^i = \begin{bmatrix} 1 \\ f^i \end{bmatrix} ; \quad \tilde{w} = \begin{bmatrix} b \\ w \end{bmatrix}$$

$$\Rightarrow \ell(\tilde{w}) = \sum_{i=1}^n (\tilde{f}^{i\top} \tilde{w} - y^i)^2$$

$$\Rightarrow \left[\sum_{i=1}^n \tilde{f}^i \tilde{f}^{i\top} \right] \tilde{w} = \sum_{i=1}^n \tilde{f}^i y^i$$

$$\boxed{\tilde{F} \tilde{F}^\top \tilde{w} = \tilde{F} y}$$

- Recap

4.1 f
A properly designed feature (or set of features) for linear regression provides a good nonlinear fit in the original feature space and, simultaneously, a good linear fit in the transformed feature space.

∴ just because we can visualize a low-dimensional regression dataset does not mean we can easily design a proper feature transformation of "by eye".

∴ it might not be obvious how to formulate a proper feature function to recover the original data generating function $y(x)$.

∴ when $x \in \mathbb{R}^d$, $d > 2 \Rightarrow$ it is even harder to "see"!

Example :

