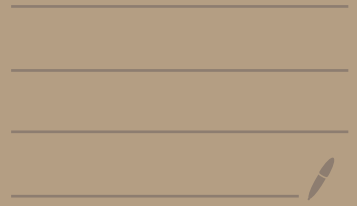


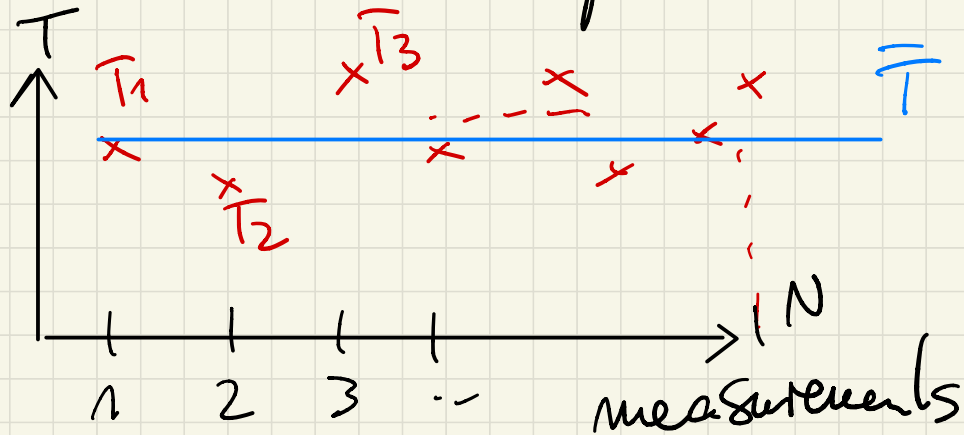
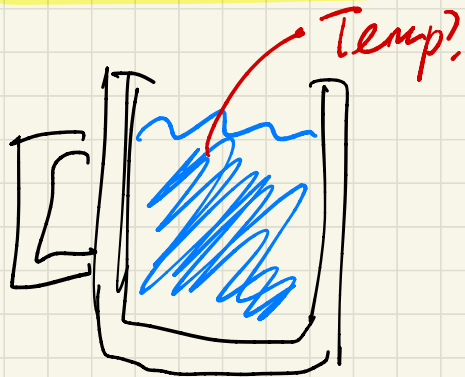
Lecture 5:

Unit 2

Optimization meets
linear algebra &
calculus 2



Example: measure the temperature



$$\bar{T} = \frac{1}{N} \sum_{i=1}^N T_i \leftarrow \text{input}$$

↑
quantity to be estimated

Complexity?

$$O(N)$$

Where does the average come from?
Define cost function $J(T)$

$$J(T) = \sum_{i=1}^N (T_i - T)^2$$

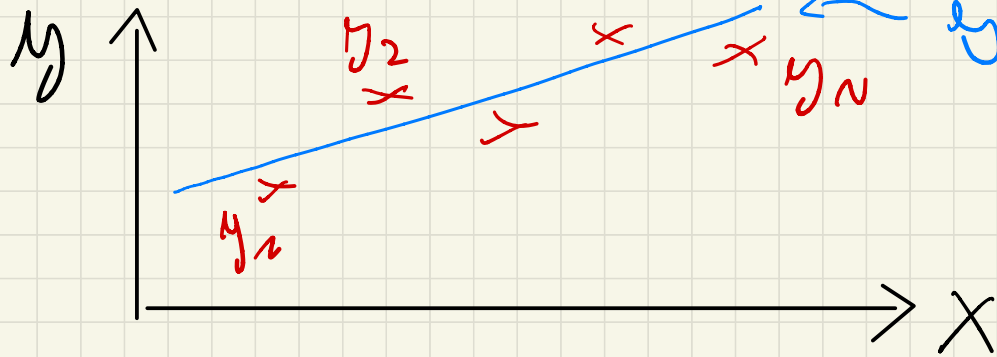
least squares

Then: $\bar{T} = \arg \min_T J(T)$

Indeed: $\frac{dJ}{dT} = 0 \iff -2 \sum_{i=1}^N T_i - T = 0$

$$\iff T = \frac{1}{N} \sum_{i=1}^N T_i \quad \square$$

Another example: linear regression



$$\bar{y} = a\bar{x} + \bar{b}$$

coefficients to be found.

Least squares:

$$\bar{a}, \bar{b} = \arg \min_{a, b} \sum_{i=1}^N \underbrace{(y_i - (ax_i + b))^2}_{J(a, b)}$$

\bar{a}, \bar{b} must satisfy: $\frac{\partial J}{\partial a} \equiv 0$, $\frac{\partial J}{\partial b} \equiv 0 \Leftrightarrow \bar{a}, \bar{b}$ as in lab 6.

What do these problems have in common?

general form: $\vec{y} = H\vec{\theta} + \vec{\epsilon}$

"noise"
(unknown)

observation operator

(model \rightarrow measurements)

model parameters
(ex: $T, \{a, b\}$)

Ex1: $H = 1$

Ex2: $H = \begin{bmatrix} 1 & x_0 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix}$

$$H \in \mathbb{R}^{m \times p}$$

$$\vec{y} \in \mathbb{R}^m$$

$$\vec{\theta} \in \mathbb{R}^p$$

Linear least squares:

$$\vec{\theta} = \arg \min_{\vec{\theta}} J(\vec{\theta}), \quad J(\vec{\theta}) = \|\vec{y} - H\vec{\theta}\|_{2}^2$$

Question: Is $\vec{\theta}$ unique?

2-norm

Stationary points: $\nabla_{\vec{\theta}} J = \vec{0}$

$$\begin{aligned} \partial_{\theta_k} J(\vec{\theta}) &= \partial_{\theta_k} \sum_i (y_i - \sum_{j=1}^n H_{ij} \theta_j)^2 \\ &= \sum_i 2 (y_i - \sum_{j=1}^n H_{ij} \theta_j) H_{ik} \end{aligned}$$

Grouping all terms $\partial_{\theta_1} \dots \partial_{\theta_p}$:

$$\nabla_{\theta} J = -2H^T (\vec{y} + H\vec{\theta}) \equiv \vec{0}$$

$$\Leftrightarrow (H^T H) \vec{\theta} = H^T \vec{y}$$

$\vec{\theta}$ is unique when $H^T H$ is invertible, H is full rank

This occurs when $p \leq n$, and
at least p rows of $H \in \mathbb{R}^{n \times p}$
are l.i.

$\uparrow \Leftrightarrow$ measurements need to
be "different"

Is the stationary point a minimum?

In one variable, one checks if $\frac{d^2 J}{d\theta^2} > 0$

In more variables, the Hessian needs to be positive definite: OK if H is full-rank
 $\Rightarrow H^T H$ is pos. def.

Complexity of $\underbrace{H^T H}_{O(p^2 m)} \vec{\theta} = \underbrace{H^T \vec{0}}_{O(p m)}$ ($p \times p$ lin system)

+ Solving linear system: $O(p^3)$

For a general matrix H ,
we can solve the optimization
in $O(p^3) + O(p^2 m) + O(p m)$ ops.

What if p, m are large?
or H is not sparse?

BREAK

Gradient descent

Given $\vec{\theta}^0$, for $k = 0, \dots$, repeat

$$\vec{\theta}^{k+1} = \vec{\theta}^k - \alpha \nabla_{\vec{\theta}} J(\vec{\theta}^k), \alpha \in \mathbb{R}.$$

Important! For using GD with several variables, all " $\vec{\theta}$ " need to have same dimensions, s.t. $[\alpha] = [J]^{-1} [\vec{\theta}]^2$

Question: How to choose α to ensure convergence?

Answer: using Taylor (as in unit 1)

$$J(\vec{\theta}^{k+1}) \approx J(\vec{\theta}^k) + (\nabla_{\vec{\theta}} J(\vec{\theta}^k)) \cdot (\vec{\theta}^{k+1} - \vec{\theta}^k)$$

$$+ \frac{1}{2} (\vec{\theta}^{k+1} - \vec{\theta}^k)^T H (\vec{\theta}^{k+1} - \vec{\theta}^k)$$

dot product

Hessian

(constant if J is quadratic)

Since in gradient descent

$$\vec{\theta}^{k+1} - \vec{\theta}^k = -\alpha \nabla_{\vec{\theta}} J(\vec{\theta}^k)$$

Then, $J(\theta^{k+1}) = J(\theta^k) - \alpha \|\nabla_{\theta} J\|^2$

The grad. descent
decrease the
cost function if \downarrow

< 0

$+ \frac{\alpha^2}{2} \|\nabla_{\theta} J\|_S^2$

$\leftarrow S$ -norm

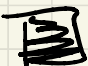
$\|x\|_S$
 $= x^T S x.$

$\Leftrightarrow \textcircled{+} \frac{\alpha^2}{2} \|\nabla_{\theta} J\|_S^2 < \alpha \|\nabla_{\theta} J\|^2$

So if: $\frac{1}{2} \max_x \frac{x^T S x}{x^T x} < \frac{1}{\alpha}, \forall x$

It holds for $\textcircled{+}$

But note that $\max_x \frac{x^T S x}{x^T x} \equiv \lambda_{\max}$
with λ_{\max} the largest eigenvalue of S .

\Rightarrow $\alpha < \frac{2}{\lambda_{\max}}$ for grad. descent
to converge. 

The closest you take $\alpha \rightarrow \frac{2}{\lambda_{\max}}$
the fastest will be the convergence,
but this also depends on the
initial guess

What about the complexity of gradient descent?

• Building $H^T H$: $O(p^2 m)$
as before

• Computations of λ_{\max} .

With power iterations (lab 7)

$$O(m_p p^2)$$

of power iterations

• Grad. Descent iters

$$O(p^2 \maxit)$$

matrix-vector multiplication

Total grad desc: $O(p^2 \cdot n) + O(n_p \cdot p^2) + O(pn \cdot n_{\text{exit}})$

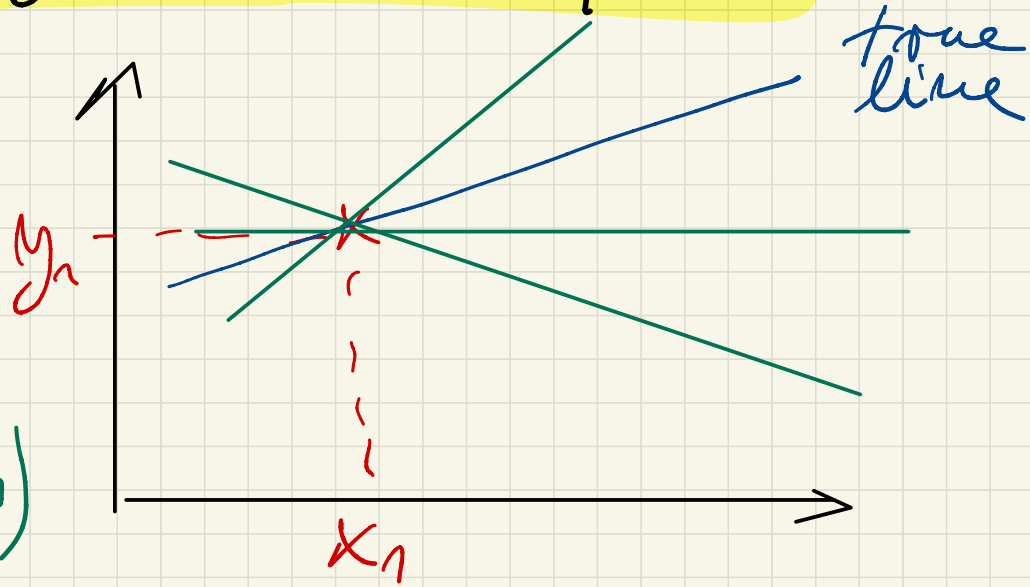
Direct ^{v/s} computation $O(p^3) + \dots$

BUT: if H is sparse then complexity computations will change
→ project.

Regularized least squares

Example:

infinite
lines passing
through (x_1, y_1)



⇒ non-unique solution

Why? 1 data point, but 2 params
to estimate (slope + intercept).

Usually two
"fixes"

Assume something
on missing data

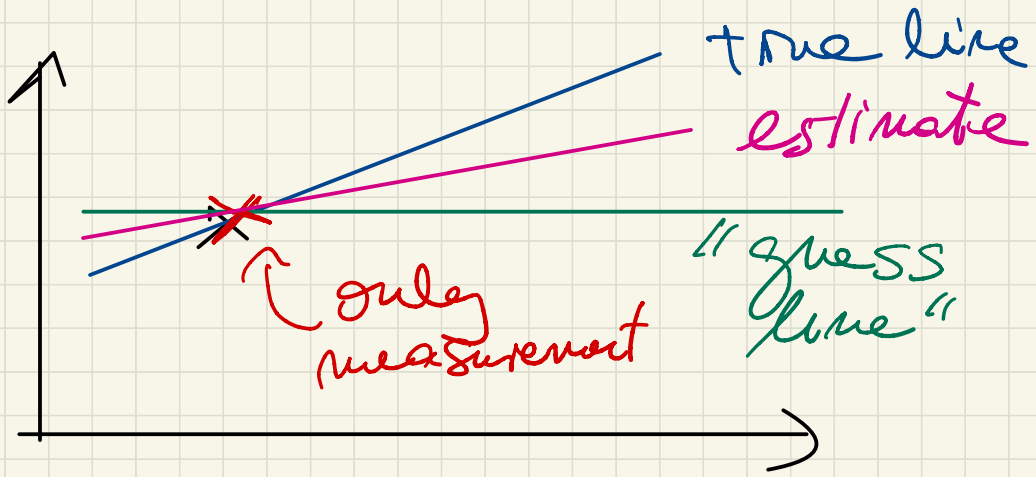
→ "create" more
measurements

→ often not good
to do.

Assume something
on parameters

↑ we will do
this

Graphically:



Mathematically:

$$J(\vec{\theta}) = \|\vec{y} - H\vec{\theta}\|^2 + \beta \|\vec{\theta} - \vec{\theta}_0\|_M^2$$

M pos def. given

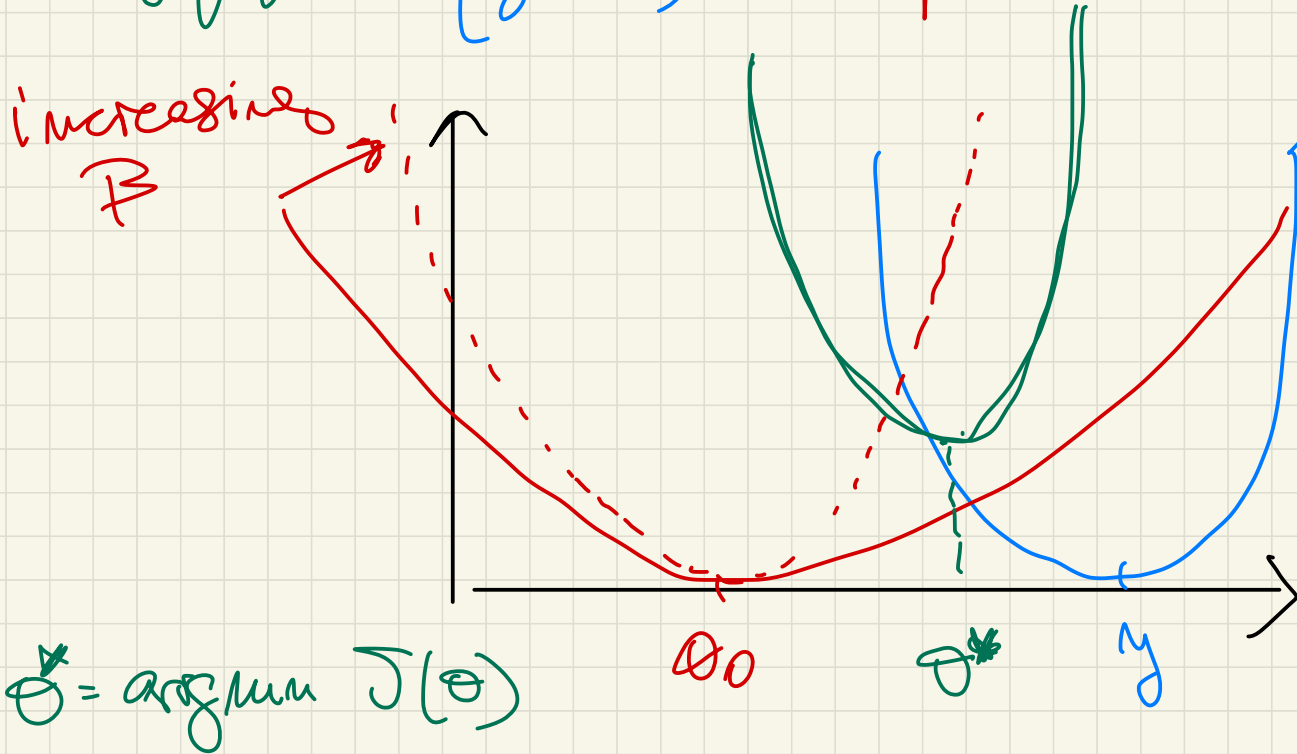
$\beta > 0$ -
given

Some
convenient
matrix norm

"prior"
known value
for the guess

Example for 1 parameter.

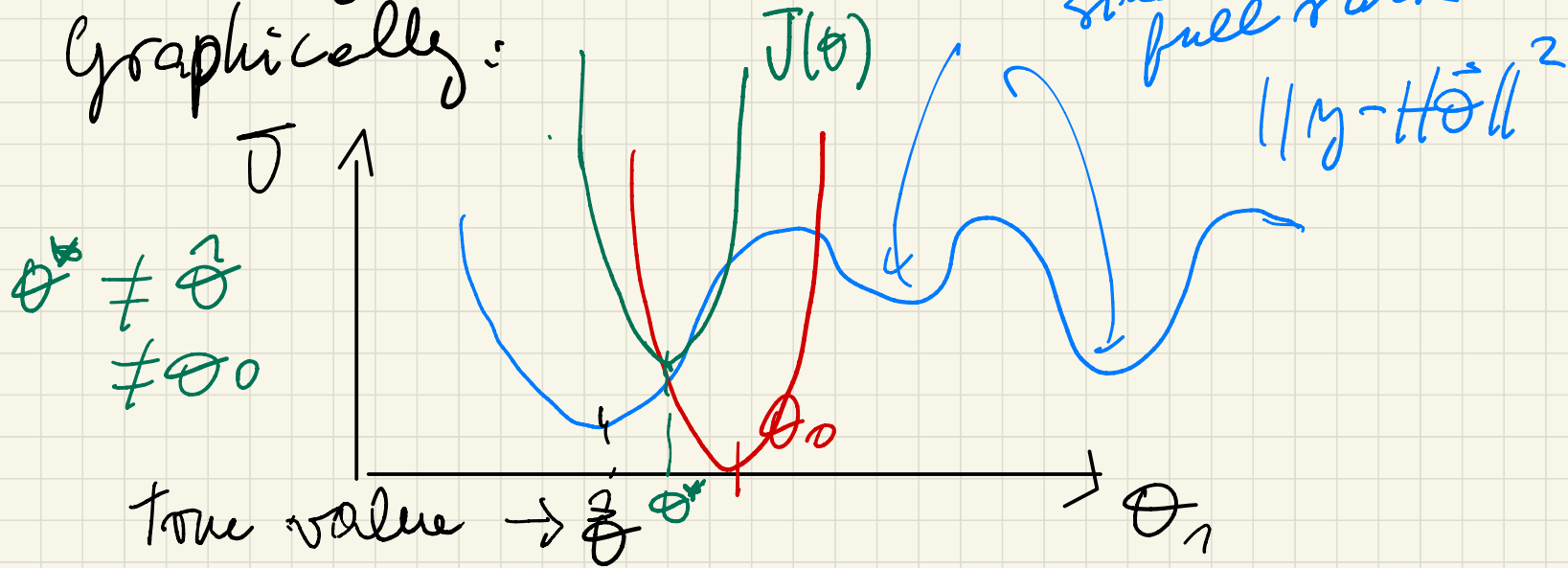
$$J(\theta) = (y - \theta)^2 + \beta (\theta - \theta_0)^2$$



Term $\beta \|\vec{\theta} - \vec{\theta}_0\|_M^2$ is called
 "regularization"

local minimum
 since H is not
 full rank

Graphically:



Remark: you also can add regularizers
 when H is full rank. (poor)

Last question: is $\vec{\theta}^*$ unique?

Yes:
$$\nabla_{\vec{\theta}} J = -2H^T \vec{y} + 2H^T H \vec{\theta} + 2\beta M (\vec{\theta} - \vec{\theta}_0)$$

So, the stationary points satisfy

$$(H^T H + \beta M) \vec{\theta}^* = H^T \vec{y} + \beta M \vec{\theta}_0$$

semi-pos
def.
(not invertible)

Sym
pos def.
 \Rightarrow invertible

Symm. pos. def. \rightarrow invertible! ~~III~~