02/22/18. Regularization.

$$y = Ax_0 + e \quad , \quad x = \text{argmin } \|Ax - y\|$$
$$= (A^T A)^{-1} A^T y \quad (\text{overdet})$$

$$\frac{\|x - x_0\|}{\|x_0\|} \leq K(A) \frac{\|e\|}{\|Ax_0\|} \quad \text{with } K(A) = \|A^+\| \|A\|$$

ex: polyn. interpolat$^{\circ}$

$$y_i = \sum c_n t_i^n + e_i$$

ex. image deblurring , fig 4.1 p 84 Bertero
$$y = h * x_0 + e \qquad (h = \text{blur kernel, PSF})$$
$$\hat{y}(k) = \hat{h}(k)\,\hat{x}_0(k) + \hat{e}(k)$$
$$\hat{x}(k) = \hat{y}/\hat{h}(k) = \hat{x}_0(k) + \hat{e}(k)/\hat{h}(k)$$

$\frac{1}{\hat{h}(k)}$ small values of $\hat{h}(k)$

$\longrightarrow e^{-k^2 \sigma^2/2}$ when spatial width $\sigma$

like a heat equation backwards.

Tikhonov regularization: push $\|x\|$ down.

$(T_\lambda)$ $\quad \min_{x} \|Ax - y\|^2 + \lambda \|x\|^2.$

$$\nabla = 0 \Rightarrow 2A^T(Ax - y) + 2\lambda x = 0$$
$$(A^T A + \lambda I)\,x = A^T y$$
$$x = (A^T A + \lambda I)^{-1} A^T y \qquad \text{unreg. when } \lambda = 0$$

In terms of the SVD:

$$A = U\Sigma V^* \qquad \text{with } U^*U = I \; ; \; V^*V = VV^* = I$$
$$A^+ = (A^TA)^{-1}A^T = (V\Sigma U^* U\Sigma V^*)^{-1} V\Sigma U^* \qquad \text{so } V^* = V^{-1}$$
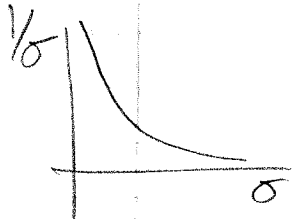$$= (V\Sigma^2 V^*)^{-1} V\Sigma U^*$$
$$= V\Sigma^{-2} V^* V\Sigma U^*$$
$$= V\Sigma^{-1} U^*$$

$$x_{LS} = A^+ y = \sum_i v_i \frac{1}{\sigma_i} u_i^* y \qquad \rightarrow \text{Fourier analysis if } A \text{ is transl-inv.}$$
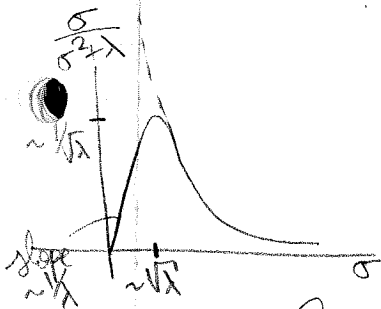
$$(A^TA + \lambda I)^{-1} A^T = (V\Sigma U^* U\Sigma V^* + I)^{-1} V\Sigma U^*$$
$$= (V\Sigma^2 V^* + VV^*)^{-1} V\Sigma U^*$$
$$= (V(\Sigma^2 + \lambda I) V^*)^{-1} V\Sigma U^*$$
$$= V(\Sigma^2 + \lambda I)^{-1} V^* V\Sigma U^*$$
$$= V(\Sigma^2 + \lambda I)^{-1} \Sigma U^*$$

$$x_{T_\lambda} = A^{+,\lambda} y = \sum_i v_i \frac{\sigma_i}{\sigma_i^2 + \lambda} u_i^* y$$

(Can do a truncated SVD instead).

Rmk: Tikhonov also proposed    (early 1960s)

$$\min_x \|Ax - y\|^2 + \lambda \|Bx\|^2 \qquad \text{Prior:}$$
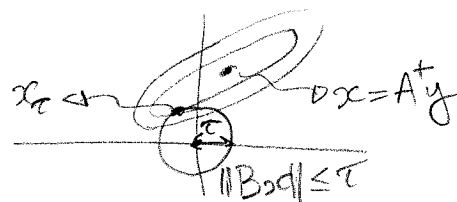
$$\text{with } B = \begin{bmatrix} 1 & -1 & & \\ & -1 & 1 & \\ & & \ddots & \\ & & -1 & 1 \end{bmatrix} \qquad \|Bx\| \text{ is small}$$

$$\Rightarrow x = (A^TA + \lambda B^TB)^{-1} A^T y$$

→ encodes smoothness of the object.
(Think $x_i = f(t_i)$ for some smooth function $f(t)$ )

ex. Bertero p. 114. Intro bias-variance.
(Hubble space telescope, 1990)

$$x_\tau \leftarrow \quad Dx = A^+ y$$
$$\|Bx\| \leq \tau$$

Equivalent formulation (constrained).

$(P_\tau)$    $\min \|Ax - y\|$   s.t.   $\|Bx\| \leq \tau$

General:    $(T_\lambda)$   $\min \boxed{f_0(x)} + \lambda \boxed{f_1(x)}$

         $(P_\tau)$   $\min f_0(x) : f_1(x) \leq \tau$

Assume   $\lambda, \tau > 0$
      $f_0, f_1$ differentiable, convex.
      Slater: $\forall \tau, \exists x : f_1(x) < \tau$.

Prop. $\bullet \forall \tau, \exists \lambda(\tau) :$ a solution $x_\tau$ of $(P_\tau)$
            is also a solution of $(T_\lambda)$
   $\bullet \forall \lambda, \exists \tau(\lambda) :$ a solution $x_\lambda$ of $(T_\lambda)$
            is also a solution of $(P_\tau)$

Pf. $\bullet$ Let $x_\tau$ be a solution to $(P_\tau)$.
    $d(x, \lambda) = f_0(x) + \lambda(f_1(x) - \tau)$.
    Slater $\Rightarrow$ strong duality, KKT holds for some $\lambda$:
   $\begin{cases} ① \nabla f_0(x_\tau) + \lambda \nabla f_1(x_\tau) = 0 \,, & ② \lambda \geq 0 \\ ③ f_1(x_\tau) - \tau \leq 0 \,, & ④ \lambda(f_1(x_\tau) - \tau) = 0. \end{cases}$

     $d$ convex in $x$, $①\Rightarrow x_\tau$ is a minimizer
     of $d$, or $f_0(x) + \lambda f_1(x)$.
   $\bullet$ Let $x_\lambda$ be a solution to $(T_\lambda)$
     $f_0(x) + \lambda f_1(x)$ convex in $x \Rightarrow$
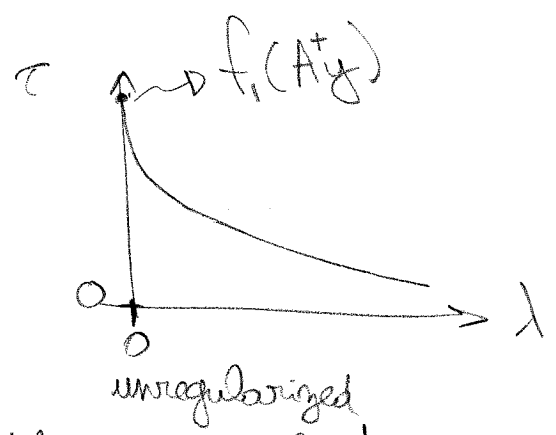     $\nabla f_0(x_\lambda) + \lambda \nabla f_1(x_\lambda) = 0$     $(1)$
     Pick $\tau = f_1(x_\lambda)$, then $(2),(3),(4)$ hold
     for $d(x, \lambda) = f_0(x) + \lambda(f_1(x) - \tau)$
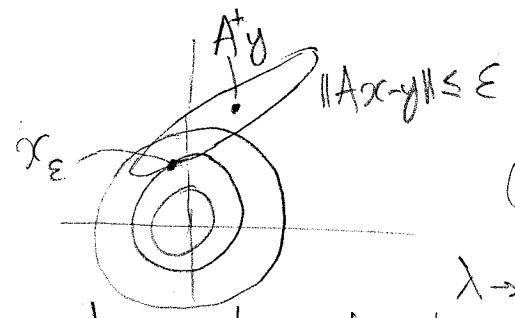     $\Rightarrow x_\lambda$ is optimal for $(P_\tau)$     $\square$

$$\tau \quad \rightarrow f_1(A^+ y)$$

$$\lambda \rightarrow \tau = f_1(A^{+,\lambda} y)$$
(other way is less direct).

unregularized

Also equivalent to

$(Q_\varepsilon)$  min $\|Bx\|$  s.t. $\|Ax-y\| \le \varepsilon$

$A^+ y$

$\|Ax - y\| \le \varepsilon$

$x_\varepsilon$

(similar argument)

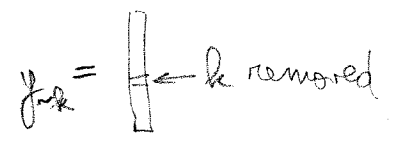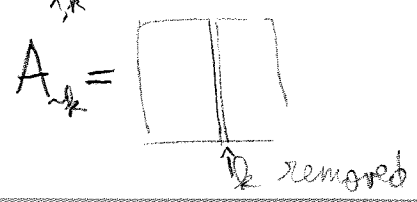$$\lambda \rightarrow \varepsilon = f_o(A^{+,\lambda} y)$$

In practice: how to choose $\tau / \lambda / \varepsilon$ ?

- provide all to user!
- estimate noise-level $\varepsilon$ such that
  $\|Ax-y\| \approx \varepsilon$
  ex. median absolute deviation (MAD)
  $$\frac{\text{median}(X_1, \dots X_m)}{0.6745}$$ on a small patch where the mean is removed.

  (Then find $\lambda / \tau$ to match $\varepsilon$)

- cross-validation
  "leave one datum out":
  $$\min_{x_{\lambda,k}} \|A_{-k} x - y_{-k}\|^2 + \lambda \|Bx\|^2 \rightarrow x_{\lambda,k}$$

  $A_{-k} = $ [matrix with column removed] $\quad$ $y_{-k} = $ [vector $\leftarrow k$ removed]

  $k$ removed

then measure prediction of remaining
component: cross validation function to min over $\lambda$:

$$CV(\lambda) = \sum_k |(A x_{\lambda,k})_k - y_k|^2$$

Hard to compute. But ($\$$)

$$CV(\lambda) = \sum_k \frac{|(A x_\lambda)_k - y_k|^2}{|1 - P_{kk}(\lambda)|^2}$$

with $\begin{cases} P(\lambda) = A(A^*A + \lambda I)^{-1} A^* \\ \quad \text{(regularized projector onto Ran A)} \\ x_\lambda = \text{argmin} \, \|Ax - y\|^2 + \lambda \|By\|^2 \end{cases}$

Issue. $V(\lambda)$ depends on the ordering of
the rows of $A$ / elements of $y$.
To fix this, consider the generalized CV
function

$$GCV(\lambda) = \sum_k \frac{|(A x_\lambda)_k - y_k|^2}{|m - \text{tr}(P(\lambda))|^2} \cdot$$

$\qquad A \in \mathbb{R}^{m \times m}$
$\qquad P(\lambda) \in \mathbb{R}^{m \times m}$

$$= \frac{\|A x_\lambda - y\|^2}{(\text{tr}(I - P(\lambda)))^2}$$

Mueller-Siltanen examples