

# Machine Learning & Inverse Problems

## Inverse Problems

$$y = Af + w$$



$$\begin{aligned} A^\top y &= (A^\top A)f + A^\top w \\ \stackrel{\text{def.}}{=} u &\stackrel{\text{def.}}{=} C \stackrel{\text{def.}}{=} r \end{aligned}$$

Regularized inversion:

$$\min_f \frac{1}{2} \|Af - y\|^2 + \lambda \|f\|^2$$

$$f_\lambda = (C + \lambda \text{Id}_p)^{-1} u$$


Exact covariance  $C$

Deterministic bounded noise  $r$

Noise level  $\varepsilon \stackrel{\text{def.}}{=} \|r\|$

## Statistical Learning

$$y = Xf + \varepsilon$$



$$\begin{aligned} \frac{1}{n} X^\top y &= \frac{1}{n} (X^\top X) f + \frac{1}{n} X^\top \varepsilon \\ \stackrel{\text{def.}}{=} u_n &\stackrel{\text{def.}}{=} C_n \stackrel{\text{def.}}{=} r_n \\ \downarrow n \rightarrow +\infty &\quad \downarrow (x_i, y_i)_i \text{ i.i.d.} \\ u &= \mathbb{E}(yx) \quad C = \mathbb{E}(xx^\top) \end{aligned}$$

Empirical risk minimization:

$$\min_f \frac{1}{2n} \|Xf - y\|^2 + \lambda \|f\|^2$$


$$f_{\lambda,n} = (C_n + \lambda \text{Id}_p)^{-1} u_n$$

Exact covariance  $C$



Noisy covariance  $C_n$

Deterministic bounded noise  $r$




Random noise  $r_n$


Noise level  $\varepsilon \stackrel{\text{def.}}{=} \|r\|$



Noise level  $\|r_n\| \sim \varepsilon = n^{-\frac{1}{2}}$






Observations  $y$






Columns of  $\Phi$

$$\min_f \|y - \Phi f\|^2 + \lambda \|f\|_2^2$$



$$\min_f \|y - \Phi f\|^2 + \lambda \|f\|_1$$



Solving  $\textcolor{red}{y} \approx \textcolor{blue}{A} \textcolor{green}{x} \in \mathbb{R}^m \quad \textcolor{blue}{A} \in \mathbb{R}^{m \times n}$

Determined ( $m = n$ ):  $\textcolor{green}{x} = \textcolor{blue}{A}^{-1} \textcolor{red}{y}$

$$\begin{array}{c|c|c} \textcolor{red}{y} & = & \textcolor{blue}{A} \times \textcolor{green}{x} \\ \hline \end{array}$$

Over-determined ( $m > n$ ):  $\min_{\textcolor{green}{x}} \|\textcolor{blue}{A}\textcolor{green}{x} - \textcolor{red}{y}\|^2$

$$\textcolor{green}{x} = (\textcolor{blue}{A}^\top \textcolor{blue}{A})^{-1} \textcolor{blue}{A}^\top \textcolor{red}{y} \stackrel{\text{def.}}{=} \textcolor{blue}{A}^+ \textcolor{red}{y}$$

$$\begin{array}{c|c|c} \textcolor{red}{y} & \approx & \textcolor{blue}{A} \times \textcolor{green}{x} \\ \hline \end{array}$$

Under-determined ( $m < n$ ):  $\min_{\textcolor{green}{x}} \{\|\textcolor{green}{x}\| ; \textcolor{blue}{A}\textcolor{green}{x} = \textcolor{red}{y}\}$

$$\textcolor{green}{x} = \textcolor{blue}{A}^\top (\textcolor{blue}{A}\textcolor{blue}{A}^\top)^{-1} \textcolor{red}{y} \stackrel{\text{def.}}{=} \textcolor{blue}{A}^+ \textcolor{red}{y}$$

$$\begin{array}{c|c|c} \textcolor{red}{y} & = & \textcolor{blue}{A} \times \textcolor{green}{x} \\ \hline \end{array}$$

$A$  ill-posed and/or noise:  $\min_{\textcolor{green}{x}} \|\textcolor{blue}{A}\textcolor{green}{x} - \textcolor{red}{y}\|^2 + \lambda \|\textcolor{green}{x}\|^2$

$$\textcolor{green}{x} = (\textcolor{blue}{A}^\top \textcolor{blue}{A} + \lambda \text{Id}_n)^{-1} \textcolor{blue}{A}^\top \textcolor{red}{y} \xrightarrow{\lambda \rightarrow 0} \textcolor{blue}{A}^+ \textcolor{red}{y}$$

$$= \textcolor{blue}{A}^\top (\textcolor{blue}{A}\textcolor{blue}{A}^\top + \lambda \text{Id}_m)^{-1} \textcolor{red}{y} \quad (\text{Woodbury identity})$$

$$\begin{array}{c|c|c} \textcolor{red}{y} & \approx & \textcolor{blue}{A} \times \textcolor{green}{x} \\ \hline \end{array}$$

# Automatic Differentiation

How to compute  $\nabla \ell_{x,y}(\theta)$ ?  $\ell_{x,y}(\theta) \stackrel{\text{def.}}{=} L(f(x, \theta), y)$

Chain rule:  $\nabla \ell_{x,y}(\theta) = [\partial f(x, \theta)]^\top (\nabla L(f(x, \theta), y))$

Linear  $f(x, \theta) = \theta \times x$ :  $\partial f(x, \theta) = \theta$ .

Non-linear  $f(x, \theta)$ : painful ... but  $\ell_{x,y}$  it is just a computer program.


Computer program  $\Leftrightarrow$  directed acyclic graph  $\Leftrightarrow$  linear ordering of nodes  $(\theta_r)_r$

forward

```
function  $\ell(\theta_1, \dots, \theta_M)$ 
  for  $r = M + 1, \dots, R$ 
    |  $\theta_r = g_r(\theta_{\text{Parents}(r)})$ 
  return  $\theta_R$ 
```


backward


```
function  $\nabla \ell(\theta_1, \dots, \theta_M)$ 
   $\nabla_R \ell = 1$ 
  for  $r = R - 1, \dots, 1$ 
    |  $\nabla_r \ell = \sum_{s \in \text{Child}(r)} \partial_r g_s(\theta) \nabla_s \ell$ 
  return  $(\nabla_1 \ell, \dots, \nabla_M \ell)$ 
```



“ $\frac{\partial \ell}{\partial \theta_r} = \sum_{s \in \text{Child}(r)} \frac{\partial \ell}{\partial \theta_s} \frac{\partial \theta_s}{\partial \theta_r}$ ”


$\nabla_r \ell(\theta)$        $\nabla_s \ell(\theta)$        $\partial_r g_s(\theta)$






$\ell^1$  ball:  $K \stackrel{\text{def.}}{=} \{x \in \mathbb{R}^n ; \sum_{i=1}^n |x_i| \leq 1\}$

Bulk:  $B \stackrel{\text{def.}}{=} \{x \in \mathbb{R}^n ; \sum_{i=1}^n |x_i|^2 \leq n^{-1}\}$



small  $n$



large  $n$

**Setup:**  $\mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}$  computable in  $K$  operations.

```
def ForwardNN(A,b,Z):
    X = []
    X.append(Z)
    for r in arange(0,R):
        X.append( rhoF( A[r].dot(X[r]) + tile(b[r],[1,Z.shape[1]]) ) )
    return X
```

*Hypothesis:* elementary operations ( $a \times b, \log(a), \sqrt{a} \dots$ )  
and their derivatives cost  $O(1)$ .

**Question:** What is the complexity of computing  $\nabla \mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ ?

Finite differences:  $\nabla \mathcal{E}(\theta) \approx \frac{1}{\varepsilon} (\mathcal{E}(\theta + \varepsilon \delta_1) - \mathcal{E}(\theta), \dots, \mathcal{E}(\theta + \varepsilon \delta_1) - \mathcal{E}(\theta))$   
 $K(n+1)$  operations, intractable for large  $n$ .

*Theorem:* there is an algorithm to compute  $\nabla \mathcal{E}$  in  $O(K)$  operations.  
[Seppo Linnainmaa, 1970]


This algorithm is reverse mode  
automatic differentiation

```
def BackwardNN(A,b,X):
    gx = lcosG(X[R],Y) # initialize the gradient
    for r in arange(R-1,-1,-1):
        M = rhoG( A[r].dot(X[r]) + tile(b[r],[1,n]) ) * gx
        gx = A[r].transpose().dot(M)
        gA[r] = M.dot(X[r].transpose())
        gb[r] = MakeCol(M.sum(axis=1))
    return [gA,gb]
```



Seppo Linnainmaa

$$x_{r+1} = g_r(x_r, \theta_r) \quad \mathcal{E}(x) = L(x_{R+1}, y)$$



$$\text{Proposition: } \forall r = R, \dots, 0, \quad \nabla_{x_r} \mathcal{E} = [\partial_{x_r} g_R(x_r, \theta_r)]^\top (\nabla_{x_{r+1}} \mathcal{E})$$

$$\nabla_{\theta_r} \mathcal{E} = [\partial_{\theta_r} g_R(x_r, \theta_r)]^\top (\nabla_{x_{r+1}} \mathcal{E})$$

*Example:* deep neural network  $x_{r+1} = \rho(A_r x_r + b_r)$


$$\nabla_{x_r} \mathcal{E} = A_r^\top M_r$$

$$\forall r = R, \dots, 0, \quad \nabla_{A_r} \mathcal{E} = M_r x_r^\top \quad M_r \stackrel{\text{def.}}{=} \rho'(A_r x_r + b_r) \odot \nabla_{x_{r+1}} \mathcal{E}$$

$$\nabla_{b_r} \mathcal{E} = M_r \mathbf{1}$$

```
def ForwardNN(A,b,Z):
    X = []
    X.append(Z)
    for r in range(0,R):
        X.append( rhoF( A[r].dot(X[r]) + tile(b[r],[1,Z.shape[1]]) ) )
    return X
```

```
def BackwardNN(A,b,X):
    gx = lossG(X[R],Y) # initialize the gradient
    for r in arange(R-1,-1,-1):
        M = rhoG( A[r].dot(X[r]) + tile(b[r],[1,n]) ) * gx
        gx = A[r].transpose().dot(M)
        gA[r] = M.dot(X[r].transpose())
        gb[r] = MakeCol(M.sum(axis=1))
    return [gA,gb]
```



$$x_{r+1} = g_r(x_r) \quad g_r : \mathbb{R}^{n_r} \rightarrow \mathbb{R}^{n_{r+1}}$$


$$\partial g_r(x_r) \in \mathbb{R}^{n_{r+1} \times n_r}$$

$$\nabla g_R(x_R) = [\partial g_R(x_R)]^\top \in \mathbb{R}^{n_R \times 1}$$



$$\partial g(x) = \partial g_R(x_R) \times \partial g_{R-1}(x_{R-1}) \times \dots \times \partial g_1(x_1) \times \partial g_0(x_0)$$

Chain rule:



Forward  
 $O(n^3)$

$$\partial g(x) = ((\dots(((\frac{A_0 \times A_1}{n_0 n_1 n_2} \times A_2) \dots \times \frac{A_{R-2} \times A_{R-1}}{n_{R-2} n_{R-1} n_R}) \times A_R \frac{}{n_{R-1} n_R}) \dots \frac{}{n_1 n_2 n_3})$$

*Complexity:* (if  $n_r = 1$  for  $r = 0, \dots, R-1$ )  $(R-1)n^3 + n^2$

Backward  
 $O(n^2)$

$$\partial g(x) = A_0 \times (\frac{A_1 \times (A_2 \times \dots \times (A_{R-2} \times (A_{R-1} \times A_R))) \dots)}{n_1 n_2} \frac{}{n_{R-1} n_R}) \dots \frac{}{n_{R-2} n_{R-1}} \frac{}{n_0 n_1})$$

*Complexity:*  $Rn^2$

*Log-sum-exp:*  $\text{LSE}_\varepsilon(x) \stackrel{\text{def.}}{=} \varepsilon \log \sum_i e^{x_i/\varepsilon}$

*Soft-max:*  $\text{SM}_\varepsilon(x) \stackrel{\text{def.}}{=} \nabla_x \text{LSE}(x) = \frac{1}{\sum_i e^{x_i/\varepsilon}} \left( e^{x_i/\varepsilon} \right)_i$



*Prop.*  $\text{LSE}_\varepsilon(x) \xrightarrow{\varepsilon \rightarrow 0} \max(x)$   $\text{SM}_\varepsilon(x) \xrightarrow{\varepsilon \rightarrow 0} \delta_{\arg\max(x)}$

LSE trick:

$$\text{LSE}_\varepsilon(x) = \text{LSE}_\varepsilon(x - \max(x)) + \max(x)$$

Unstable

Stable



Soft-max:  $\text{SM}(u)_k \stackrel{\text{def.}}{=} \frac{e^{u_k}}{\sum_\ell e^{u_\ell}}$



Log-sum-exp:  $\text{LSE}(u) \stackrel{\text{def.}}{=} \log \sum_k e^{u_k}$

Logit model:  $\mathbb{P}(\text{Class}(x) = k) \stackrel{\text{def.}}{=} \text{SM}((\langle x, w_\ell \rangle)_\ell)_k$

Training data:  $(x_i, y_i)_i$   $y_{i,k} = \mathbb{P}(\text{Class}(x_i) = k)$

Logistic classification:

$$\min_{(w_k)_k} \sum_i \text{LSE}((\langle x_i, w_k \rangle)_k) - \sum_{i,k} y_{i,k} \langle x_i, w_k \rangle$$




$\ell^p$  “norms”

$$\|x\|_p^p \stackrel{\text{def.}}{=} \sum_i |x_i|^p$$


Lasso / Basis-Pursuit:

$$\min_x \|x\|_1 + \frac{1}{2\lambda} \|Ax - y\|^2$$
$$\min_{Ax=y} \|x\|_1$$


$\xleftarrow[\lambda \rightarrow 0]$




Robert  
Tibshirani




David Donoho




$p = 0.5$



$p = 0.75$



$p = 1$




Non-convex

Non-sparse

$$\text{Elastic net: } x_\lambda \in \operatorname{argmin}_x \frac{1}{2\lambda} \|Ax - y\|^2 + (1 - \theta)\|x\|_1 + \frac{\theta}{2}\|x\|_2^2$$

Regularization path:  $\lambda \longmapsto x_\lambda$



*Supervised learning:* Observations:  $(a_i, y_i)_i$ , parametric model:  $g(x, a)$

Regression:	$y_i \approx g(x, a_i)$	$\ell(y, y') =  y - y' ^2$
-------------	-------------------------	----------------------------

Classification:	$y_i \approx \theta(g(x, a_i))$	$\ell(y, y') = \log(1 + e^{-yy'})$
	$\theta(u) = (1 + e^u)^{-1}$	

Empirical risk minimization:	$\min_x f(x) = \frac{1}{n} \sum_i \ell(g(x, a_i), y_i)$
------------------------------	---

$$\min_x f(x)$$

