## Inverse Problems:

$$\hat{x} = f_\theta(y) \quad \text{(Inference function)}$$



$x$ (unknown) (to be inferred) → THE WORLD → $y$ (data) (measure-able) → INVERSION / INFERENCE → $\hat{x}$ (estimate)

$\Theta$ (parameters)

*All sensor problems are inverse problems

---

$$\hat{x} = \boxed{B} \; \sigma \left\{ \boxed{A} \; \boxed{y} + \boxed{b} \right\}$$

(function applied) ← $\sigma$

col vector   offset vector

---

## Supervised Machine Learning:



Training data → $x_k$ / $(x_k, y_k)$ pairs / $y_k$ (data) → ML Training done via minimizing the Loss → $\Theta$

---

## Single Layer Dense NN:



$y \to \boxed{Ay} \to \boxed{+} \to \sigma(\cdot) \to \boxed{B} \to \hat{x}$

$b$

$$\hat{x} = B\sigma(Ay + b)$$

$$\begin{cases} A \in \mathbb{R}^{N_1 \times N_y} \\ b \in \mathbb{R}^{N_1} \\ \sigma : \mathbb{R}^{N_1} \to \mathbb{R}^{N_1} \\ B \in \mathbb{R}^{N_x \times N_1} \end{cases}$$

matrix of multiplicative weights
column vector of additive offsets
point-wise activation function
matrix of multiplicative weights

---

## Abstraction:

$$\hat{x} = f_\theta(y) = B\sigma(Ay + b)$$



$y \to \boxed{Ay} \to + \to \sigma(\cdot) \to \boxed{B} \to \hat{x}$

$b$

$$\Theta = (A, B, b) \quad \text{("the set of all NN params.")}$$

### Single Layer NN Flow Diagram



$(N_y = 3)$  (input layer)   (hidden layer)   $(N_x = 2)$  (output layer)
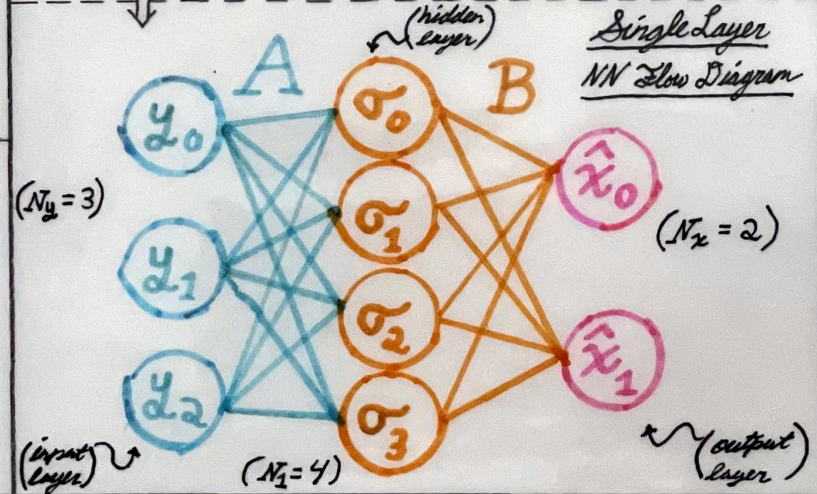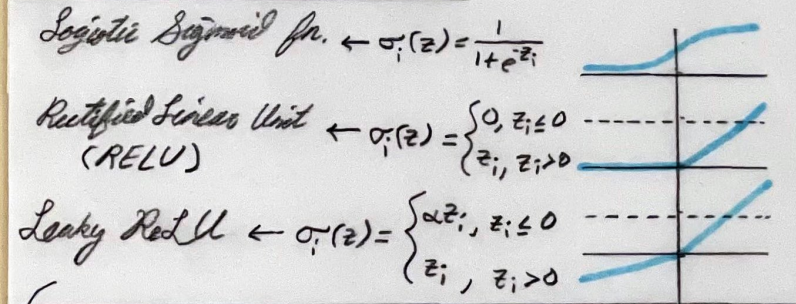
$(N_1 = 4)$

---

## Gradient matrix of Softmax Function:

*Dense matrix → mostly non-zero values

$$[\nabla\sigma(z)]_{i,j} = \frac{1}{\sum_k e^{z_k}} \left( e^{z_i} \delta_{ij} - \frac{e^{z_i} e^{z_j}}{\sum_k e^{z_k}} \right)$$
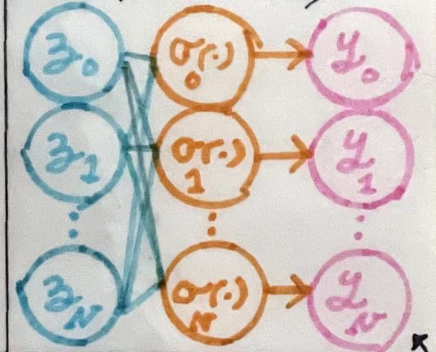
*Slow 2 compute

$$e^{z_i} e^{z_j} =$$

$$\sigma_i(z) = \frac{e^{z_i}}{\sum_j e^{z_j}} \leftarrow \left( \frac{\text{SOFTMAX}}{\text{ACT. FN}} \right)$$

*Joint act. fn. can be interpreted as a probability density



$z_0 \to \sigma(\cdot)_0 \to y_0$
$z_1 \to \sigma(\cdot)_1 \to y_1$
$z_N \to \sigma(\cdot)_N \to y_N$

---

## Point-Wise Activation Functions:

Logistic Sigmoid fn. $\leftarrow \sigma_i(z) = \dfrac{1}{1 + e^{-z_i}}$

Rectified Linear Unit (RELU) $\leftarrow \sigma_i(z) = \begin{cases} 0, & z_i \leq 0 \\ z_i, & z_i > 0 \end{cases}$

Leaky ReLU $\leftarrow \sigma_i(z) = \begin{cases} \alpha z_i, & z_i \leq 0 \\ z_i, & z_i > 0 \end{cases}$

---

$\sigma : \mathbb{R}^N \to \mathbb{R}^N$ (point-wise act. fn)



$z_0 \to \sigma_0(\cdot) \to y_0$
$z_1 \to \sigma_1(\cdot) \to y_1$
$z_N \to \sigma_N(\cdot) \to y_N$

$z \to \boxed{\sigma(\cdot)} \to y$

("is a point-wise activation function")

---

Gradient Matrix of a Point-wise Act. fn. is
• diagonal
• sparse (many zeros)
• fast 2 compute

$$\nabla\sigma(z) = \left[ \frac{\partial \sigma_i(z)}{\partial z_j} \right] = \begin{bmatrix} d_1 & & 0 \\ & \ddots & \\ 0 & & d_N \end{bmatrix}$$

$(N \times N)$

$$\left( \text{where } d_i = \frac{\partial \sigma(z_i)}{\partial z_i} \right)$$

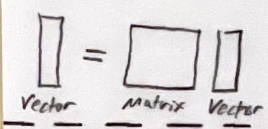Standard Encoding $\leftarrow \hat{x} \in \{0, \ldots, M-1\}$ (each val repr. a class)

1-Hot Encoding (for classification) $\leftarrow \hat{x} \in \mathbb{R}^M$ s.t. $\hat{x}_i = \begin{cases} 1, & \text{class} == i \\ 0, & \text{class} \neq i \end{cases}$

---

## M-dimensional Simplex:

$$\mathscr{S}^M = \left\{ x \in \mathbb{R}^M \,\middle|\, \forall i, x_i \geq 0 \;\&\; 1 = \sum_{i=0}^{M-1} x_i \right\}$$

$$\hat{x} \in \mathscr{S}^M \subset \mathbb{R}^M \quad \left( \text{like a probability density per class} \right)$$

*Continuous function on a convex set
*Allows for representation of prob. densities



(3D simplex)   *Good for optimization

# Tensors

(from "differential geometry", used by Einstein when formulating the theory of General Relativity)

**Def.: Tensor**
- The generalization of a
  { matrix operator, for dim >2
  { vector data to >1 dim

$$\text{Vector} \quad = \quad \text{Matrix} \quad \text{Vector}$$

**MATRICES ARE OPERATORS**

OUT ⟸ [box] ⟸ IN  ($j_2$, $i$, $j_1$)

Example of a tensor generalizing a matrix operator for >2 dims

($j_2$, $j_1$) or ($j_2$, $j_1$)

Example of a tensor generalizing a vector to a higher dimensional space / representation

---

$$f_\theta(y) = B\sigma(Ay+b)$$

$$y \to \boxed{Ay} \to \oplus \to \boxed{\sigma(\cdot)} \to \boxed{B} \to \hat{x}$$

$A \qquad b \qquad B$

$$\delta y \to \boxed{\nabla_y f} \to \oplus \to \oplus \to \oplus \to \delta x$$

$\boxed{\nabla_A f} \quad \boxed{\nabla_b f} \quad \boxed{\nabla_B f}$

$\delta A \qquad \delta b \qquad \delta B$

**Parameters:** $\theta = (A, b, B)$

**Gradients** (at the input) $\left\{ \begin{array}{l} \nabla_\theta f_\theta(y) \\ \nabla_y f_\theta(y) \text{ or } \nabla_y f \end{array} \right.$

---

**Partial Derivative of** $A_{ji}$ **w.r.t.** $y_i$

$$\frac{\partial [Ay]_j}{\partial y_i} \xrightarrow{\text{(only depends on)}} \frac{\partial (A_{ji} y_i)}{\partial y_i} = A_j^i$$

$$\left( i.e., \; [\nabla_y(Ay)]_j^i = A_j^i \right)$$

($j$th input)

$= \left( \begin{array}{c} j\text{th} \\ \text{output} \end{array} \right) \to \boxed{A} \; \boxed{y} = Ay$

**Delta Function:**

$$\delta_j^i = \left\{ \begin{array}{ll} 1, & i=j \\ 0, & i \neq j \end{array} \right.$$

$$\Downarrow$$

$$G_j^n = G_i^k \delta_j^i \xrightarrow{\text{(Sum over the i's)}} \left( \begin{array}{c} \text{Gradient} \\ \text{w.r.t.} \\ \text{Vector} \end{array} \right) \nabla_y(Ay) = A$$

( Gradient w.r.t. Matrix )
$$[\nabla_A(Ay)]_{j_1, j_2}^i = \delta_{j_1}^i y^{j_2}$$

tensor notation $\to [\nabla_y(Ay)]_j^i = A_j^i$

---

(Contravariant ⟺ Data) col vectors

(Covariant ⟺ Operator) row vectors

*Note: differentiating a 1D obj w.r.t. a 2D obj yields a ⟹ 3D obj

---

**Contravariant vectors:** $x = gy$

*Column vectors representing data describe the position of something
* $y^j$ for $0 \leq j < N$ & $x \in \mathbb{R}$

**Covariant vector:**
* Row vector (gradient vec) that operates on data
* $g_j$ for $0 \leq j < J$

**Einstein Notation:**
$$x = g_j y^j = \sum_{j=0}^{N-1} g_j y^j$$
{ - leave out Summation
{ - Always sum over any 2 indices that occur twice

---

(scalar)
$$\boxed{x} = \boxed{g} \qquad \boxed{y}$$
$(1 \times N)$

(Covariant) ⤳
(Contravariant) $(N \times 1)$

$$x = Gy$$

* 1D Contravariant Vector
$y^j$ for $0 \leq j \leq N_y$
$x^j$ for $0 \leq j \leq N_x$

* 2D Tensor (i.e., matrix)
$G_j^i$ for $0 \leq i < N_x$ & $0 \leq j < N_y$

---

**Vector-Matrix Products:**
$$\boxed{x} = \boxed{G} \quad \boxed{y}$$
$(N \times 1) \qquad (N \times N) \qquad (N \times 1)$

(Covariant rows) (Contravariant)

$$x^i = G_j^i y^j = \sum_{j=0}^{N_y - 1} G_j^i y^j$$

---

**Tensor Products:** (sum over)

$$x^{i_1, i_2} = G_{j_1, j_2}^{i_1, i_2} y^{j_1, j_2}$$

**2D Contravariant Vectors:**
$y^{j_1, j_2}$ for $0 \leq j_1, j_2 < N_y$
$x^{i_1, i_2}$ for $0 \leq i_1, i_2 < N_x$

**4D Tensor:**
$G_{j_1, j_2}^{i_1, i_2}$ for $0 \leq i_1, i_2 < N_x$
& for $0 \leq j_1, j_2 < N_y$

---

Example: $x^i = G_{j_1, j_2}^i y^{j_1, j_2}$

$$\boxed{x} = \boxed{G} \quad \boxed{y}$$
(3D tensor) (2D tensor)

G is a tensor with 2D covariant input & 2D contravariant output

$G \in \mathbb{R}^{N_x \times N_y}$ ← General Linear Transform

# Gradient Descent Optimization (GD)

## GD Algorithm:

$$d = -\nabla \mathcal{L}(\theta) \leftarrow 1 \times N_x \text{ (row vector)}$$

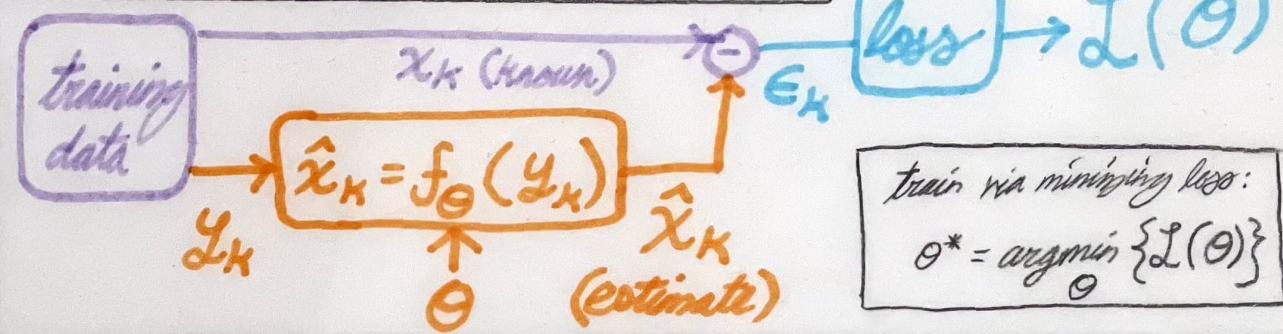repeat until converged {
$$d \leftarrow -\nabla \mathcal{L}(\theta)$$
$$\theta \leftarrow \theta + \alpha \, d^t$$
}

If $\alpha$ is too $\begin{cases} \text{small} \rightarrow \text{SLOW} \\ \text{large} \rightarrow \text{UNSTABLE} \end{cases}$

## Steepest Descent: (EXPENSIVE)

* compute best $\alpha$ via __line search__

repeat until converged {
$$d \leftarrow -\nabla \mathcal{L}(\theta)$$
$$\alpha^* \leftarrow \arg\min_{\alpha} \{ L(\theta + \alpha d^t) \}$$
$$\theta \leftarrow \theta + \alpha^* d^t$$
}

## Coordinate Descent:

* Update 1 param. at a time
* Fast but requires many updates

## Slow Convergence of Gradient Descent

* Sensitive to __condition number__ (CN) of the problem, no perfect step size choice

* Solution → Newton's Method, correct for local 2nd derivative ... "sphere the ellipse" (too computational / difficult)

* Alternative methods such as momentum



$$CN = \frac{\sigma_{max}}{\sigma_{min}} = 2 \qquad CN = \frac{\sigma_{max}}{\sigma_{min}} = 10$$

---

* Compute gradient via chain rule
* Compute Adjoint gradient
* Back Propagate



train via minimizing loss:
$$\theta^* = \arg\min_{\theta} \{ \mathcal{L}(\theta) \}$$

$$d = -\nabla \mathcal{L}(\theta) \quad \text{(gradient is a row vector)} \quad (1 \times N_x)$$

$$\Rightarrow \begin{array}{l} \text{Gradient Descent Convergence} \\ \text{do while } \xi \\ \quad d \leftarrow -\nabla \mathcal{L}(\theta) \\ \quad \theta \leftarrow \theta + \alpha d^t \, r_{(\text{rate})} \\ \xi \text{ converges} \end{array}$$

* The Norm of a vector is the square root of the sum of the square of its elements (euclidean distance)

---

## Loss Gradient

*Note that the "Adjoint Matrix" / Conjugate transpose / Hermitian transpose is just the transpose for real matrices → $A^H$ or $A^* = A^T$

$$\nabla_\theta \mathcal{L}_{MSE}(\theta) = \nabla_\theta \left\{ \frac{1}{K} \sum_{k=0}^{K-1} \| x_n - f_\theta(y_k) \|^2 \right\} = \frac{1}{K} \sum_{k=0}^{K-1} \nabla_\theta \left\{ \| x_k - f_\theta(y_k) \|^2 \right\}$$

$$= \frac{2}{K} \sum_{k=0}^{K-1} (x_k - f_\theta(y_k))^t \, \nabla_\theta (x_k - f_\theta(y_k)) = \frac{-2}{K} \sum_{k=0}^{K-1} (x_k - f_\theta(y_k))^t \, \nabla_\theta f_\theta(y_k)$$
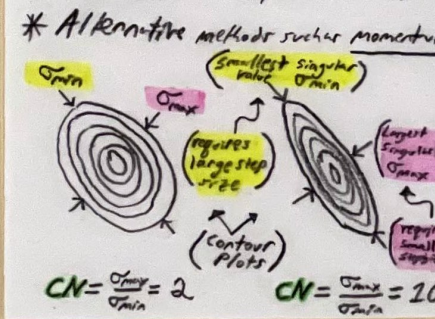
$$\therefore \quad -\nabla_\theta \mathcal{L}_{MSE}(\theta) = \frac{2}{K} \sum_{k=0}^{K-1} \underbrace{(x_k - f_\theta(y_k))^t}_{} \underbrace{\nabla_\theta f_\theta(y_k)}_{}$$
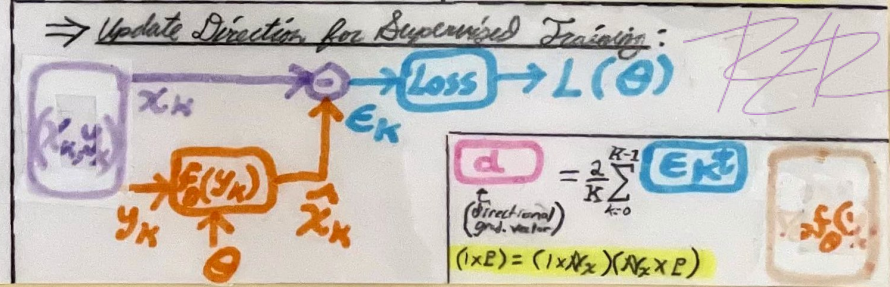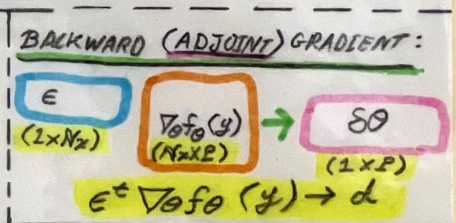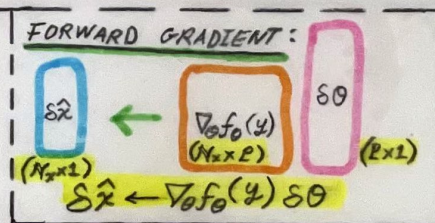
(Summation over training data.)  (Prediction Error)  (Gradient of Function)

(Gradient of inference function is enabled by AUTOMATIC DIFFERENTIATION)

(ie, Back-propagation for NN)

---

### Error Vector $(1 \times N_x)$
$$= \epsilon_k^t = (x_k - f_\theta(y_k))^t$$
$$(\because x_k = f_\theta(y_k) + error)$$

### Param. Vector $(P \times 1)$
$$= [-\nabla_\theta \mathcal{L}_{MSE}]^t = \begin{pmatrix} \text{dimensionality of param. vector} \end{pmatrix}$$

$j \longrightarrow$ (param. dim)

(output dim) $i$

#### Inference Function Gradient
$$\frac{\partial [f_\theta(y_k)]_i}{\partial \theta_j} \quad (N_x \times P) = \nabla_\theta f_\theta(y_k) \quad \begin{pmatrix} \text{gradient of} \\ \text{function} \end{pmatrix}$$

---

### FORWARD GRADIENT:



$$\delta \hat{x} \leftarrow \nabla_\theta f_\theta(y) \, \delta\theta$$

$\delta\hat{x}$ $(N_x \times 1)$ , $\nabla_\theta f_\theta(y)$ $(N_x \times P)$ , $\delta\theta$ $(P \times 1)$

### BACKWARD (ADJOINT) GRADIENT:



$$\epsilon^t \nabla_\theta f_\theta(y) \rightarrow d$$

$\epsilon$ $(1 \times N_x)$ , $\nabla_\theta f_\theta(y)$ $(N_x \times P)$ , $\delta\theta$ $(1 \times P)$

---

$\Rightarrow$ ## Update Direction for Supervised Training:



$$d = \frac{2}{K} \sum_{k=0}^{K-1} \epsilon_k^t \quad \begin{pmatrix} \text{directional} \\ \text{grad. vector} \end{pmatrix}$$

$$(1 \times P) = (1 \times N_x)(N_x \times P)$$

# Backprop. 4 CNNs

**Adj. Grad. 4 Convolution** | **Adjoint Gradient** | **Gradient of Convolution** (output w.r.t. input)

---

\* 2 functions required per _Node_

1) **forward propagation:**



$$x \leftarrow f(Z, \theta)$$

\* Gradient of output w.r.t. _weights_

$$x_i = z_i * w_i = \sum_j z_{i-j} w_j \quad \Leftrightarrow \quad x = Aw$$

$$\frac{\partial x_i}{\partial w_j} = A_{i,j} = z_{i-j}$$

$$[A^t]_{i,j} = A_{j,i} = z_{j-i}$$

$$\delta_i = \sum_j z_{j-i} \epsilon_j$$

( Auto correlating $\epsilon_j$ w/ time reverse of $y_j$ )

$$[A^t]_{i,j} = A_{j,i} = W_{j-i}$$

$$\delta_i = \sum_j W_{j-i} \epsilon_j = W_{-i} * \epsilon_i$$



\* Adjoint uses the time-reversed impulse response

( TRANSPOSE OF A )

\* fwd prop., then find gradient of the output w.r.t. the input } A



$$x_i = W_i * z_i = \sum_j W_{i-j} z_j \quad \Leftrightarrow \quad x = Az$$

$$\frac{\partial x_i}{\partial z_j} = A_{i,j} = W_{i-j} \quad \text{(TOEPLITZ MATRIX)}$$

---

2) **adjoint gradient**



( multiplying $\epsilon$ by the adjoint gradient )

$$g = (g_\omega, g_b)$$

$$[\delta, g_\omega, g_b] \leftarrow G(\epsilon, Z, \theta)$$



$$g_w \quad g_b$$

\* REMINDER → SW implementation of the convolution is really a correlation operation

$$f(y) = \sigma\left( W_{(j_1, j_2), i}^{j_3} * Z_i^{(j_1, j_2)} + b^{j_3} \right)$$

$$\begin{array}{l} A = \nabla_Z f_\theta(z) \leftarrow \text{Gradient w.r.t. input, } Z \\ B_w = \nabla_w f_\theta(z) \leftarrow \text{Gradient w.r.t. filter weights, } W \\ B_b = \nabla_b f_\theta(z) \leftarrow \text{Gradient w.r.t. offsets, } b \end{array}$$

\* Convolutions are commutative i.e., $\omega_i * z_i \Leftrightarrow z_i * \omega_i$

---

Single Layer CNN Example:



$$f_\theta(z) = \sigma(W * z + b)$$
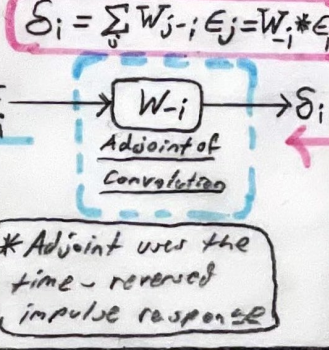
$$\nabla_\theta f_\theta(z) = \left[ \nabla_w f_{(\omega, b)}(z) \quad \nabla_b f_{(\omega, b)}(z) \right]^t$$

Adj. Grad. w.r.t. $\theta = (\omega, b)$, & then get $\nabla_z f_\theta(z)$ ↰ (gradient w.r.t. input)

(Adjoint Gradient) (Adjoint function Gradient) (Error Vector)  Backpropagating the error, $\epsilon$

$$\delta = A_{j,i} = \frac{\partial [f_\theta(y_x)]_j}{\partial z_i} \cdot \epsilon$$

↳ (Huge Matrix)

$$(N_y \times 2) \qquad (N_y \times N_x) \qquad (N_x \times 1)$$

$\zeta$ (Input Gradients)

## FAST ADJOINT GRADIENT APPROACH

⇒ _Directly compute_ the outputs w/o constructing the gradient matrix _prior_ & requiring both memory & computational resources ("fast" ∵ A is never computed!)