

Good: 1 label y_i **drawn** $\sim x_i^2$ Predict: $w_i^* = \frac{y_i}{x_i}$

$$\mathbb{E}_{i}[L(w_{i}^{*})] = 2L(w^{*})$$

$$\mathbb{E}_{i}[w_{i}^{*}] = \sum_{i} \frac{x_{i}^{2}}{\|\mathbf{x}\|^{2}} \frac{y_{i}}{x_{i}} = w^{*}$$

 $P(i) \quad w_i^*$

Unbiased estimates for linear regression viavolumesampling

Michał Dereziński

Subsampling for linear regression

Given: *n* points $\mathbf{x}_i \in \mathbb{R}^d$ with hidden labels $y_i \in \mathbb{R}$ **Goal**: Minimize loss $L(\mathbf{w}) = \sum_{i} (\mathbf{x}_{i}^{\top} \mathbf{w} - y_{i})^{2}$ over all *n* points

Select
$$S = \{4, 6, 9\}$$

Receive y_4, y_6, y_9

 \mathbf{X}

Simple strategy: Solve the subproblem, $\mathbf{w}^*(S) = \mathbf{X}_S^{+\top} \mathbf{y}_S$

Volume sampling

 $S \subseteq \{1..n\}$ chosen w.p.

 \sim squared volume of parallelepiped spanned by the $\{\mathbf{x}_i : i \in S\}$

Distribution over all *d*-element subsets *S*:

 $P(S) = \det(\mathbf{X}_S \mathbf{X}_S^{\top}) / Z$

Normalization factor obtained via Cauchy-Binet formula:

$$Z = \sum_{S:|S|=d} \det(\mathbf{X}_S \mathbf{X}_S^{\top}) =$$

Loss expectation formula

Theorem For a volume-sampled set *S* of size *d*,

$$\mathbb{E}\left[L(\mathbf{w}^*(S))\right] = (d+1) \ L(\underbrace{\mathbf{w}^*}_{\mathbb{E}[\mathbf{w}^*(S)]}),$$

if **X** is in general position

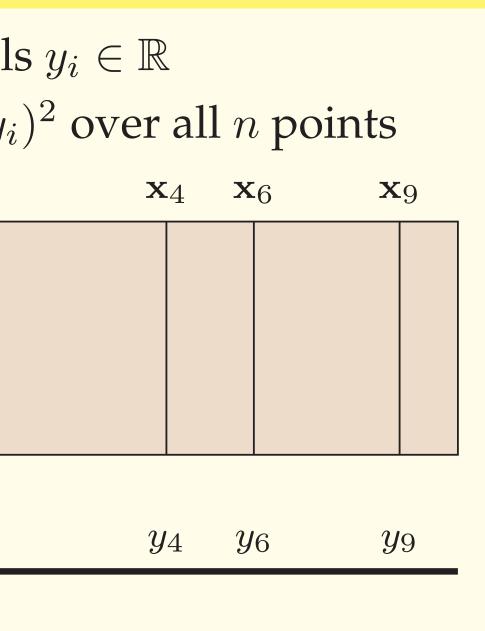
- distribution does not depend on labels - no range restrictions!

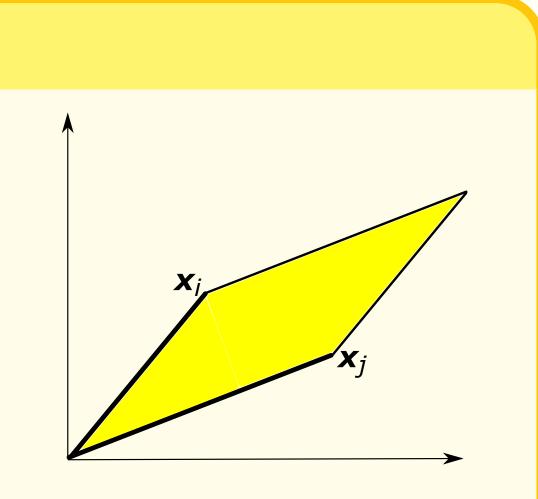
- x_{max} (furthest from 0) is bad - any deterministic choice bad



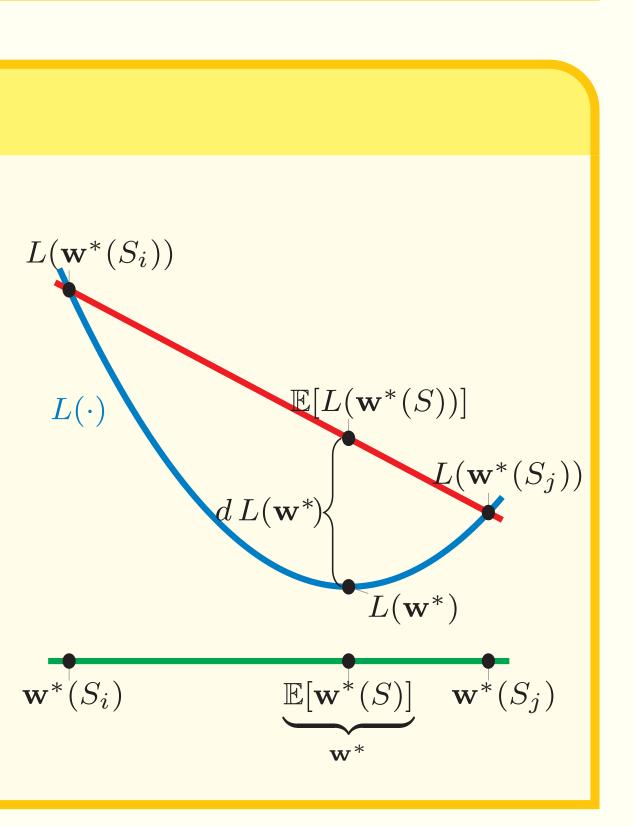
Manfred K. Warmuth

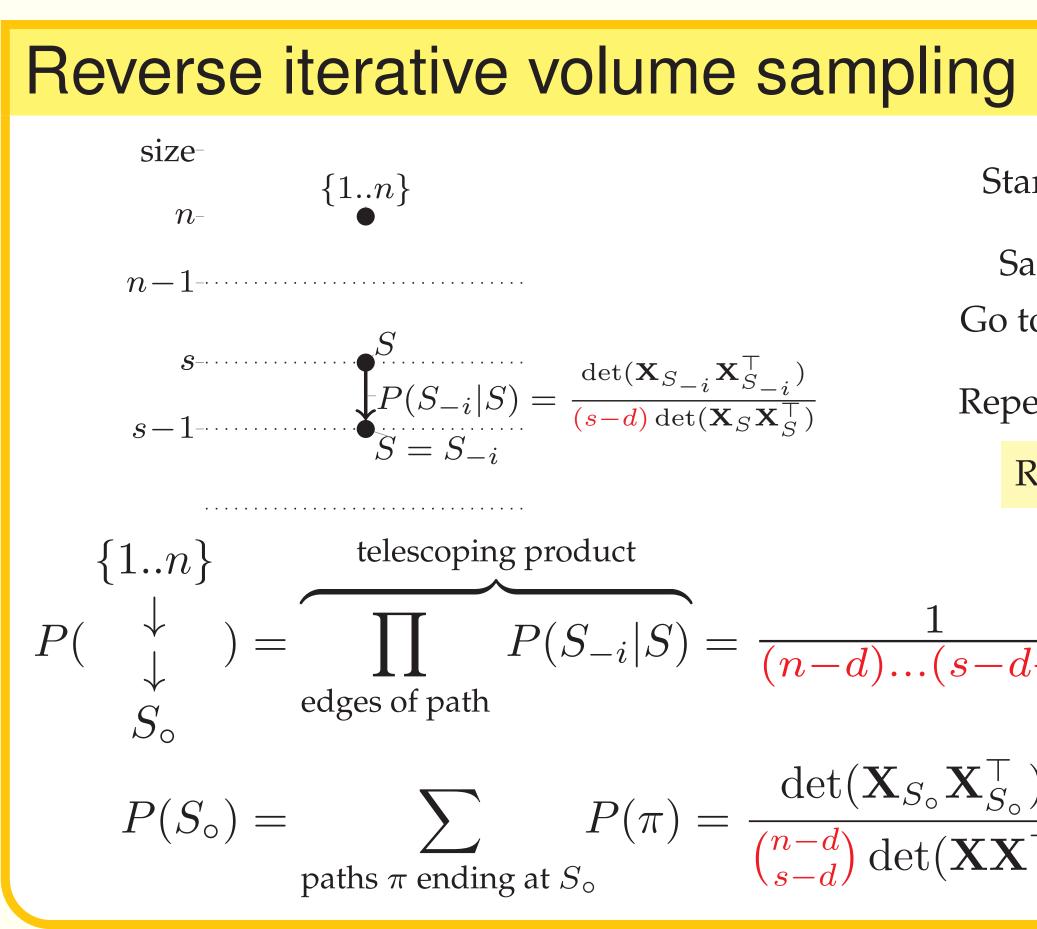






 $\det(\mathbf{X}\mathbf{X}^{\top})$





Key trick: To each subset *S* assign a formula $\mathbf{F}(S)$ st $\mathbf{F}(S) = \sum_{i \in S} P(S_{-i}|S) \mathbf{F}(S_{-i}).$

Then: $\mathbb{E}_S[\mathbf{F}(S)] = \mathbf{F}(\{1..n\})$

Expectation formulas for $(\mathbf{XI}_S)^+$ 1. $\mathbb{E}[(\mathbf{XI}_S)^+] = \mathbf{X}^+$ 2. $\mathbb{E}[(\mathbf{X}_S \mathbf{X}_S^{\top})^{-1}] = \frac{n-d+1}{s-d+1} (\mathbf{X} \mathbf{X}^{\top})^{-1}$ (variance) $((\mathbf{XI}_S)^+)^2$

Corollary: $\mathbb{E}[\mathbf{w}^*(S)] = \mathbb{E}[(\mathbf{X}\mathbf{I}_S)^{+\top}\mathbf{y}] = \mathbf{X}^{+\top}\mathbf{y} = \mathbf{w}^*$

Averaging unbiased estimators

Let $\widehat{\mathbf{y}}(S) = \mathbf{X}^{\top} \mathbf{w}^{*}(S)$. If $\mathbf{w}^{*}(S)$ is unbiased ($\mathbb{E}[\mathbf{w}^{*}(S)] = \mathbf{w}^{*}$), then:		
loss bound	variance bound	
$\mathbb{E}[L(\mathbf{w}^*(S))] \le (1+c) L(\mathbf{w}^*) \Longleftrightarrow$	$\mathbb{E}[\ \widehat{\mathbf{y}}(S) - \mathbb{E}[\widehat{\mathbf{y}}(S)]\ ^2] \le c L(\mathbf{w}^*)$	
Take average of k i.i.d. samples of size s: $\overline{\mathbf{w}}^* = \frac{1}{k} \sum_{j=1}^k \mathbf{w}^*(S_j)$,		
$\mathbb{E}\left[L\left(\overline{\mathbf{w}}^*\right)\right] \le \left(1 + \frac{c}{k}\right) \ L(\mathbf{w}^*) \Longrightarrow$	$\frac{s c}{\epsilon}$ labels achieve $(1 + \epsilon)L(\mathbf{w}^*)$	
With size <i>d</i> volume sampling, we need d^2/ϵ labels. Is d/ϵ possible?		
Open: Is there unbiased estimator with $s = O(d)$ and $c = O(1)$?		



	Start with $S = \{1n\}$
	Sample index $i \in S$ Go to set $S_{-i} = S - \{i\}$
$\frac{\det(\mathbf{X}_{S_{-i}}\mathbf{X}_{S_{-i}}^{\top})}{\mathbf{s}-d)\det(\mathbf{X}_{S}\mathbf{X}_{S}^{\top})}$	Repeat until desired size
	Runtime: $O(n^2d)$
uct	$\det(\mathbf{X}_{S_{\circ}}\mathbf{X}_{S_{\circ}}^{\top})$
$_{-i} S) = {(n-d)}.$	$\frac{1}{\dots(s-d+1)} \frac{\det(\mathbf{X}S_{\circ}\mathbf{X}S_{\circ})}{\det(\mathbf{X}\mathbf{X}^{\top})}$
$\det(\mathbf{X}$	$(\mathbf{X}_{\alpha}^{\top})$

 $P(\pi) =$

$$\frac{\det(\mathbf{A}_{S_{\circ}}\mathbf{A}_{S_{\circ}})}{\binom{n-d}{s-d}\det(\mathbf{X}\mathbf{X}^{\top})}$$

