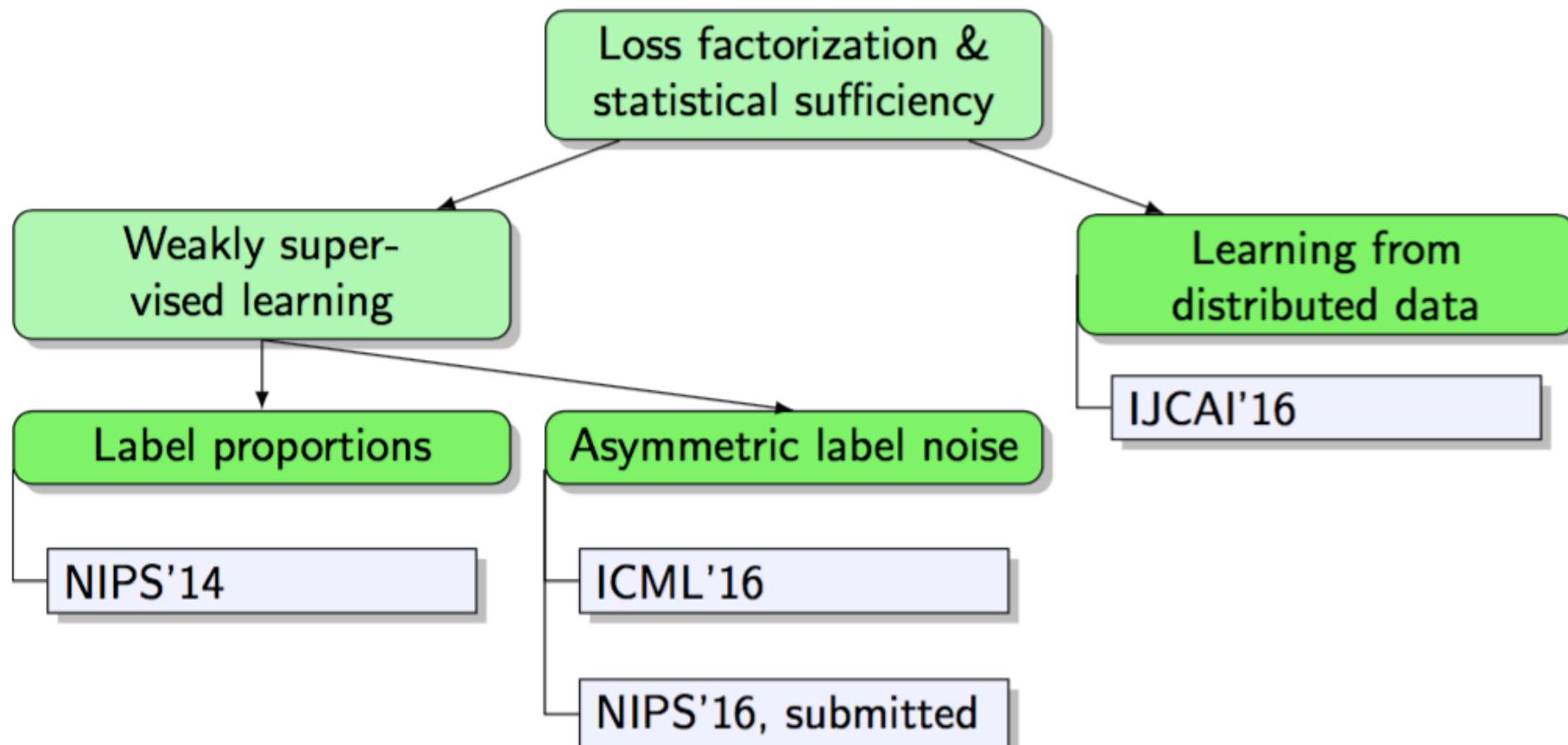


Weakly supervised learning via statistical sufficiency

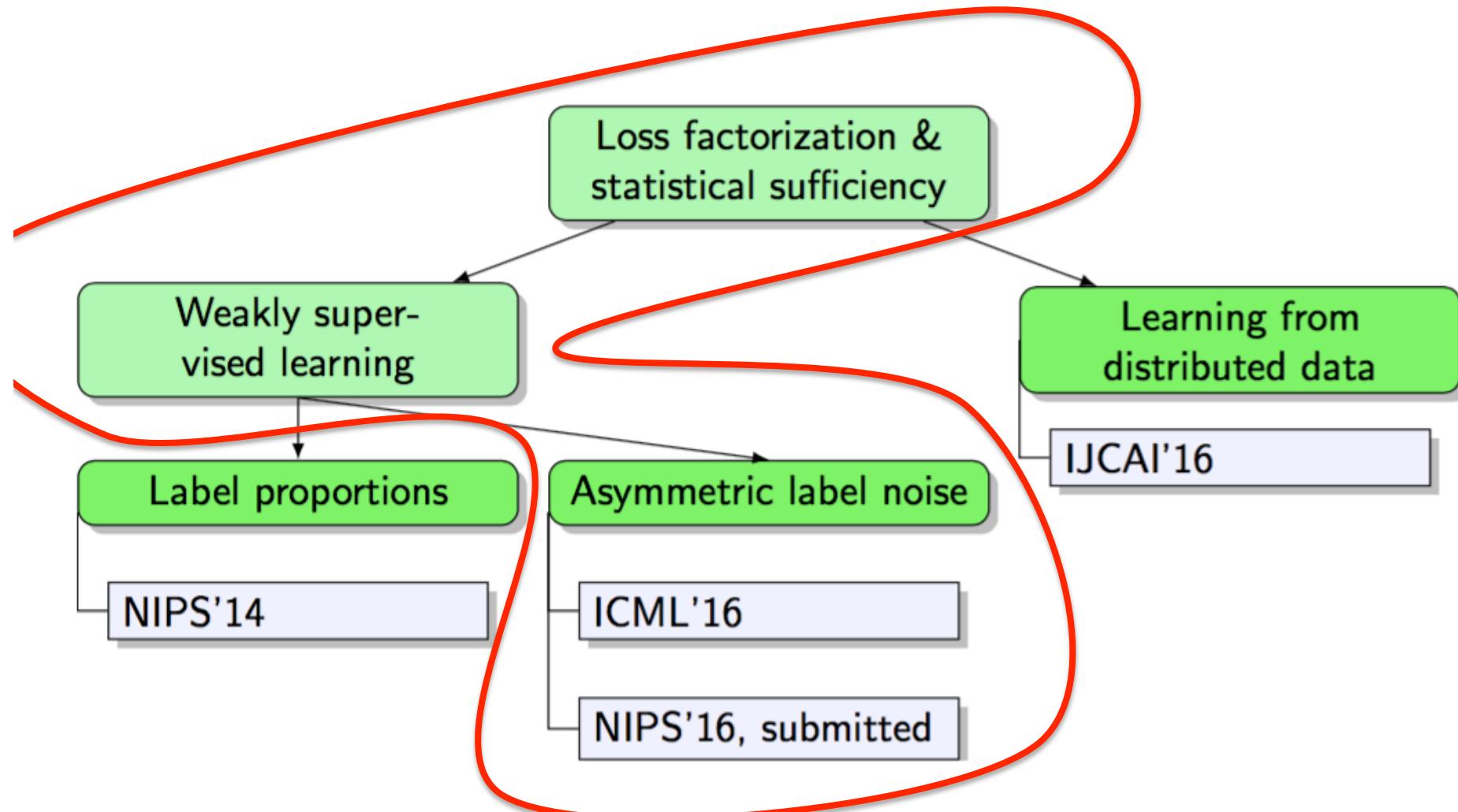
Giorgio Patrini

4 August 2016

Content organization



Content organization



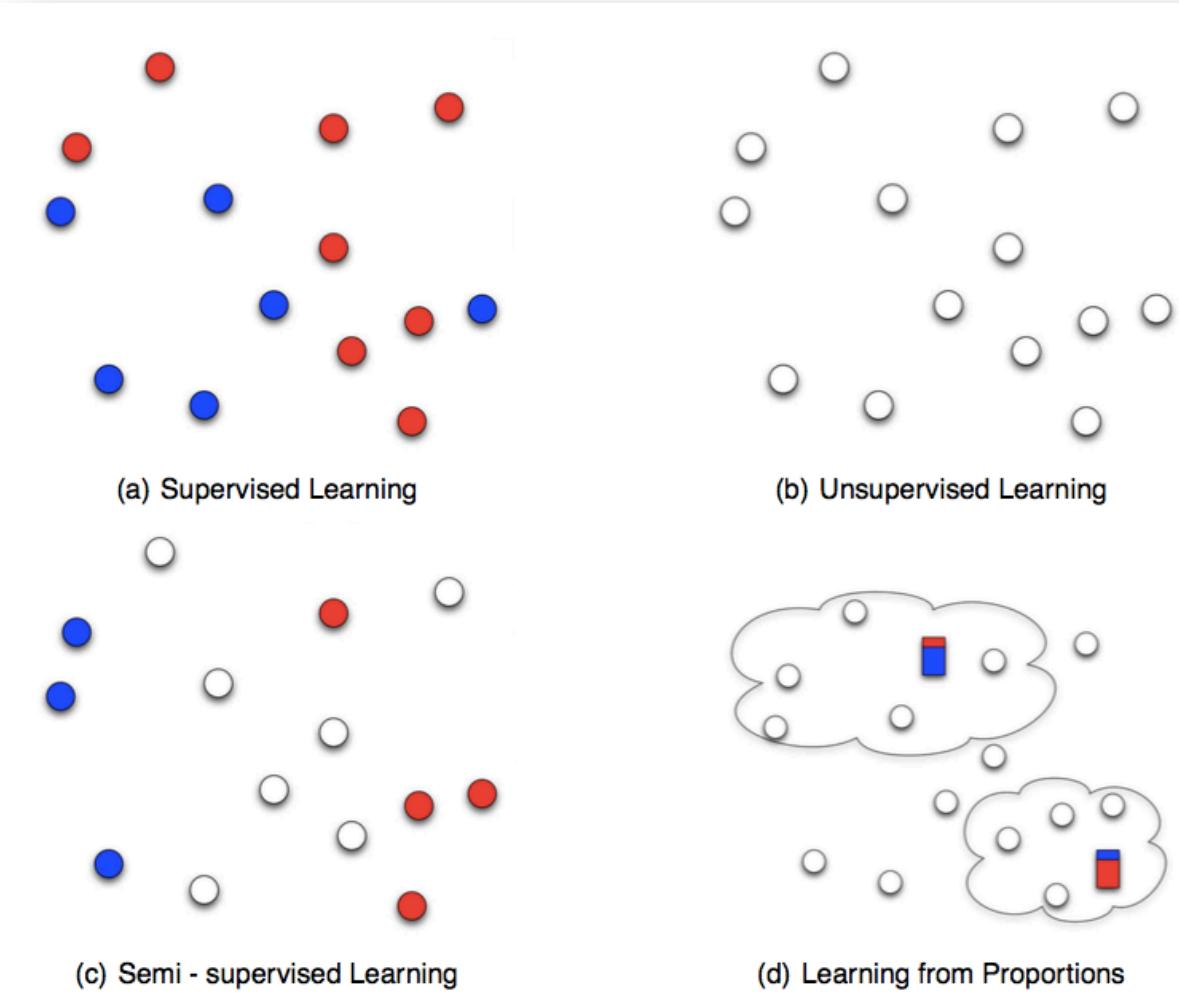
Outline

- Motivations
- Statistical sufficiency & loss factorization
- Asymmetric label noise I: theory with linear models
- Asymmetric label noise II: neural nets
- Other learning settings and conclusion

Outline

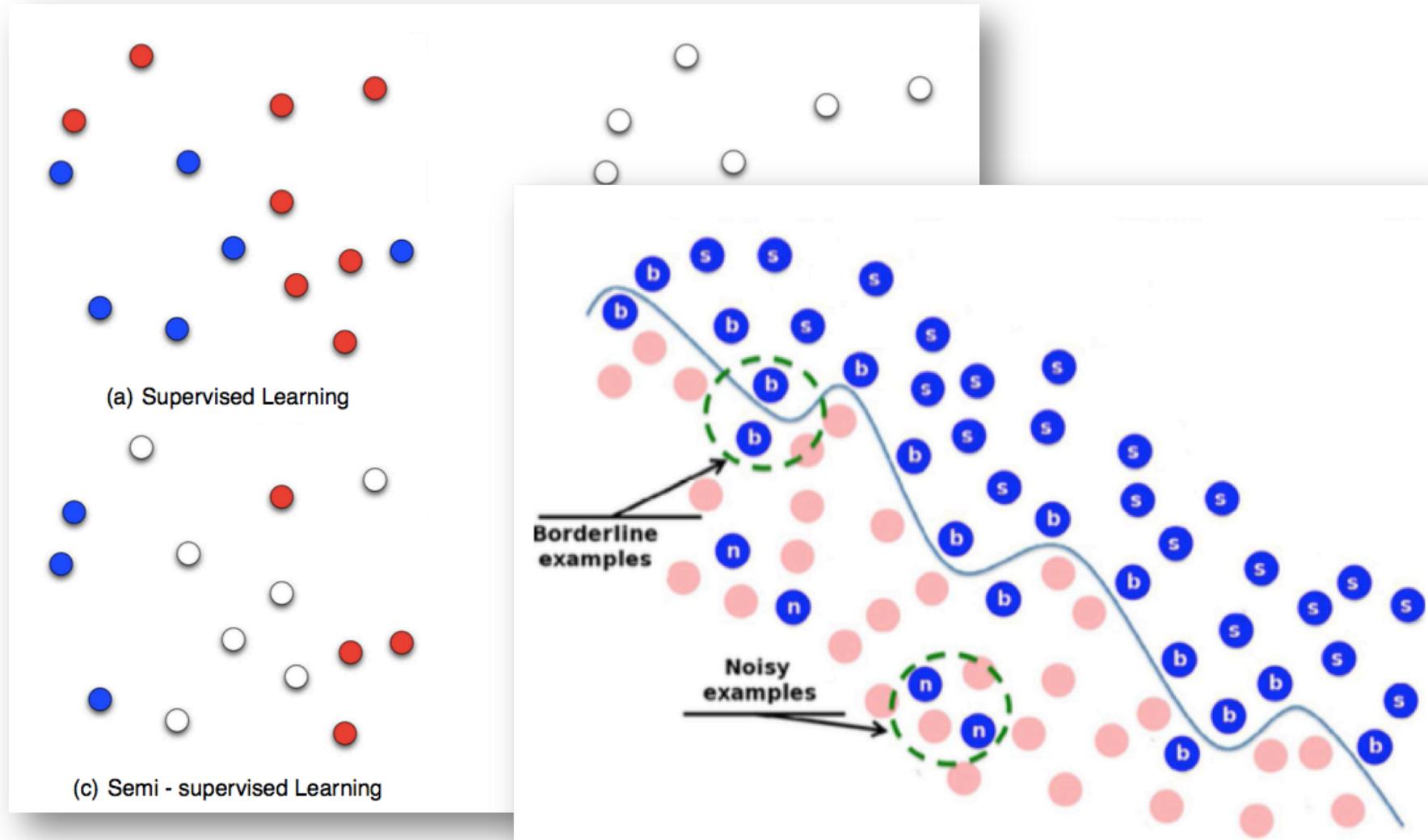
- **Motivations**
- Statistical sufficiency & loss factorization
- Asymmetric label noise I: theory with linear models
- Asymmetric label noise II: neural nets
- Other learning settings and conclusion

Motivations



[Quadrianto et al.'09]

Motivations



[sci2s.ugr.es/noisydata]

(Fully) supervised learning

- Binary classification
 $\mathcal{S} = \{(\mathbf{x}_i, y_i), i \in [m]\}$ sampled from \mathcal{D} over $\mathbb{R}^d \times \{-1, 1\}$
- Learn a linear model $h \in \mathcal{H}$
- Minimize the empirical risk associated with a surrogate loss $\ell(x)$

$$\operatorname{argmin}_{h \in \mathcal{H}} \mathbb{E}_{\mathcal{S}} [\ell(yh(\mathbf{x}))] = \operatorname{argmin}_{h \in \mathcal{H}} R_{\mathcal{S}, \ell}(h)$$

{1, ..., m}



Weakly supervised learning

$$\mathcal{D} \xrightarrow{\text{corrupt}} \tilde{\mathcal{D}} \xrightarrow{\text{sample}} \tilde{\mathcal{S}}$$

- **Weak** labels may be wrong, missing, multi-instance label constraints, etc.
- Marginal of X is the same
- The goal is unchanged: generalize to \mathcal{D}

Risk with weak labels is non-sense

$$\mathbb{E}_{\tilde{S}} [\ell(yh(\mathbf{x}))] \quad ?$$

A sample of proposed solutions

- Individual labels are corrupted: engineer a robust loss

$$\operatorname{argmin}_h \mathbb{E}_{\tilde{S}}[\tilde{\ell}(y h(\boldsymbol{x}))]$$

A sample of proposed solutions

- Individual labels are corrupted: engineer a robust loss

$$\operatorname{argmin}_h \mathbb{E}_{\tilde{S}}[\tilde{\ell}(y h(\boldsymbol{x}))]$$

- Some labels are missing: regularize with unlabelled data

$$\operatorname{argmin}_{h \in \mathcal{H}} \mathbb{E}_{\tilde{S}}[\ell(y h(\boldsymbol{x}))] + \lambda \operatorname{REG}(\{\boldsymbol{x}\}, h)$$

A sample of proposed solutions

- Individual labels are corrupted: engineer a robust loss

$$\operatorname{argmin}_h \mathbb{E}_{\tilde{S}}[\tilde{\ell}(y h(\boldsymbol{x}))]$$

- Some labels are missing: regularize with unlabelled data

$$\operatorname{argmin}_{h \in \mathcal{H}} \mathbb{E}_{\tilde{S}}[\ell(y h(\boldsymbol{x}))] + \lambda \operatorname{REG}(\{\boldsymbol{x}\}, h)$$

- Label knowledge on sets: enforce the constraints

$$\operatorname{argmin}_{h \in \mathcal{H}, y \in \{\pm 1\}} \mathbb{E}_{\tilde{S}}[\ell(y h(\boldsymbol{x}))] + \lambda \operatorname{CONSTR}(\{(\boldsymbol{x}, y)\}, h)$$

Drawbacks

A mix of

- Need to dream up new:
 - losses (possibly non-convex)
 - regularizers/constraints
 - optimization algorithms
- Label as latent variable => non-convex objective
- No unified approach

Solution principles

- Divide and conquer (decoupling)
Treat labels issues and learning separately
- Do not reinvent the wheel (modularity)
Let's reuse well-known algorithms for risk minimization
- Computational laziness
“One should solve the problem directly and never solve a more general problem as an intermediate step” [V. Vapnik'98]
- Statistical sufficiency
If we know what is **sufficient** for learning, poor labels are not a problem

A 2-step framework

- (1) Estimate a sufficient statistic μ for the labels from weak supervision

$$\tilde{\mathcal{S}} \rightarrow \mu$$

- (2) Plug it into a standard loss ℓ and call any algorithm for empirical risk minimization

$$\operatorname{argmin}_{h \in \mathcal{H}} R_{\tilde{S}, \ell}(h, \mu)$$

A unifying approach

Learning from label proportions with

- logistic loss [Quadrianto et al. '09]
- symmetric proper loss [Patrini et al. '14]

Learning with noisy labels with

- logistic loss [Gao et al. '16]

“Indirect supervision” with

- logistic loss [Raghunathan et al. '16]

Outline

- Motivations
- **Statistical sufficiency & loss factorization**
- Asymmetric label noise I: theory with linear models
- Asymmetric label noise II: neural nets
- Other learning settings and conclusion

Statistical sufficiency

- Intuition: a sufficient statistic aggregates from data all information about the model parameters

Statistical sufficiency

- Intuition: a sufficient statistic aggregates from data all information about the model parameters
- Definition: μ is sufficient for θ wrt. Y when for each pair of outcomes y, y' we have

$$\frac{P(\theta|\mu(y))}{P(\theta|\mu(y'))} \text{ does not depend on } Y \iff \mu(y) = \mu(y')$$

Exponential family \Leftrightarrow logistic loss

- Conditional exponential family

$$p(y|x) = \exp\{\langle \theta, yx \rangle - \log \sum_y \exp\langle \theta, yx \rangle\}$$

label independent

Exponential family \Leftrightarrow logistic loss

- Conditional exponential family

$$p(y|\boldsymbol{x}) = \exp\{\langle\boldsymbol{\theta}, y\boldsymbol{x}\rangle - \log \sum_y \exp\langle\boldsymbol{\theta}, y\boldsymbol{x}\rangle\}$$

label independent

- Log-likelihood leads to logistic loss

$$\sum_{i=1}^m \log \sum_y e^{y\langle\boldsymbol{\theta}, \boldsymbol{x}_i\rangle} - \langle\boldsymbol{\theta}, \boldsymbol{\mu}\rangle = \sum_{i=1}^m \log \left(1 + e^{-2y_i\langle\boldsymbol{\theta}, \boldsymbol{x}_i\rangle}\right)$$

sufficiency
of $\boldsymbol{\mu}$ for y

Mean operator & linear-odd losses

- Mean operator

$$\boldsymbol{\mu}_{\mathcal{S}} \doteq \mathbb{E}_{\mathcal{S}}[y\boldsymbol{x}] = \frac{1}{m} \sum_{i=1}^m y_i \boldsymbol{x}_i \in \mathbb{R}^d$$

Mean operator & linear-odd losses

- Mean operator

$$\boldsymbol{\mu}_{\mathcal{S}} \doteq \mathbb{E}_{\mathcal{S}}[y\boldsymbol{x}] = \frac{1}{m} \sum_{i=1}^m y_i \boldsymbol{x}_i \in \mathbb{R}^d$$

- Linear-odd loss, a -LOL

$$\exists a \in \mathbb{R}, \frac{1}{2} (\ell(x) - \ell(-x)) = \ell_o(x) = ax$$

generic x
argument

Loss factorization

- Linear model h
- Linear-odd loss $\frac{1}{2}(\ell(x) - \ell(-x)) = \ell_o(x) = ax$
- Given a sample \mathcal{S} , define a “double sample”

$$\mathcal{S}_{2x} \doteq \{(\mathbf{x}_i, \sigma), i \in [m], \sigma \in \{\pm 1\}\}$$

Loss factorization

- Linear model h
- Linear-odd loss $\frac{1}{2}(\ell(x) - \ell(-x)) = \ell_o(x) = ax$
- Given a sample \mathcal{S} , define a “double sample”

$$\mathcal{S}_{2x} \doteq \{(x_i, \sigma), i \in [m], \sigma \in \{\pm 1\}\}$$

Then:

$$R_{\mathcal{S}, \ell}(h) = \frac{1}{2} R_{\mathcal{S}_{2x}, \ell}(h) + a \cdot h(\mu_{\mathcal{S}})$$

label independent



smoothness
nor convexity
of ℓ required

Loss factorization: proof

$$R_{\mathcal{S}, \ell}(h) =$$

$$= \mathbb{E}_{\mathcal{S}} \left[\ell(yh(\mathbf{x})) \right] \quad \text{even + odd}$$

$$= \frac{1}{2} \mathbb{E}_{\mathcal{S}} \left[\ell(yh(\mathbf{x})) + \ell(-yh(\mathbf{x})) + \ell(yh(\mathbf{x})) - \ell(-yh(\mathbf{x})) \right]$$

$$= \frac{1}{2} \mathbb{E}_{\mathcal{S}_{2x}} \left[\ell(\sigma h(\mathbf{x})) \right] + \mathbb{E}_{\mathcal{S}} \left[\ell_o(h(y\mathbf{x})) \right]$$

$$= \frac{1}{2} R_{\mathcal{S}_{2x}, \ell}(h) + a \cdot h(\boldsymbol{\mu}_{\mathcal{S}})$$

sufficiency
of $\boldsymbol{\mu}$ for y

linear ℓ_o and h

Sufficiency with losses

μ is sufficient for y when $\forall \ell \in \mathcal{L}, \forall h \in \mathcal{H}$
and for each pair of learning samples $\mathcal{S}, \mathcal{S}' :$

$R_{\mathcal{S}, \ell}(h) - R_{\mathcal{S}', \ell}(h)$ does not depend on $y \iff \mu_{\mathcal{S}} = \mu_{\mathcal{S}'}$

Sufficiency with losses

μ is sufficient for y when $\forall \ell \in \mathcal{L}, \forall h \in \mathcal{H}$
and for each pair of learning samples $\mathcal{S}, \mathcal{S}'$:

$R_{\mathcal{S}, \ell}(h) - R_{\mathcal{S}', \ell}(h)$ does not depend on $y \iff \mu_{\mathcal{S}} = \mu_{\mathcal{S}'}$

Sufficiency as a consequence of factorization

$$R_{\mathcal{S}, \ell}(h) = \frac{1}{2} R_{\mathcal{S}_{2x}, \ell}(h) + a \cdot h(\mu_{\mathcal{S}})$$

Linear-odd losses: examples

	loss ℓ	odd term ℓ_o
LOL	$\ell(x)$	ax
ρ -loss	$\rho x - \rho x + 1$	$-\rho x$ ($\rho \geq 0$)
unhinged	$1 - x$	$-x$
perceptron	$\max(0, -x)$	$-x$
double-hinge	$\max(-x, 1/2 \max(0, 1 - x))$	$-x$
SPL	$a_\ell + \ell^\star(-x)/b_\ell$	$-x/(2b_\ell)$
logistic	$\log(1 + e^{-x})$	$-x/2$
square	$(1 - x)^2$	$-2x$
Matsushita	$\sqrt{1 + x^2} - x$	$-x$

Linear-odd losses: examples

	loss ℓ	odd term ℓ_o
LOL	$\ell(x)$	ax
ρ -loss	$\rho x - \rho x + 1$	$-\rho x$ ($\rho \geq 0$)
unhinged	$1 - x$	$-x$
perceptron	$\max(0, -x)$	$-x$
double-hinge	$\max(-x, 1/2 \max(0, 1 - x))$	$-x$
SPL	$a_\ell + \ell^\star(-x)/b_\ell$	$-x/(2b_\ell)$
logistic	$\log(1 + e^{-x})$	$-x/2$
square	$(1 - x)^2$	$-2x$
Matsushita	$\sqrt{1 + x^2} - x$	$-x$

- Bonus: convex LOLs are calibrated when $a < 0$

Generalization bound

- Loss is a -LOL and Lipschitz
- Bounded feature and model spaces
- Bounded loss $C \doteq \max_x \ell(x)$
- Let $\hat{\boldsymbol{\theta}} \doteq \operatorname{argmin}_{\boldsymbol{\theta} \in \mathcal{H}} R_{\mathcal{S}, \ell}(\boldsymbol{\theta})$

Generalization bound

- Loss is a -LOL and Lipschitz
- Bounded feature and model spaces
- Bounded loss $C \doteq \max_x \ell(x)$
- Let $\hat{\boldsymbol{\theta}} \doteq \operatorname{argmin}_{\boldsymbol{\theta} \in \mathcal{H}} R_{\mathcal{S}, \ell}(\boldsymbol{\theta})$

Then for any $\delta > 0$, with probability at least $1 - \delta$:

$$\begin{aligned} R_{\mathcal{D}, \ell}(\hat{\boldsymbol{\theta}}) - \inf_{\boldsymbol{\theta} \in \mathcal{H}} R_{\mathcal{D}, \ell}(\boldsymbol{\theta}) &\leq O\left(\frac{1}{\sqrt{m}}\right) + \\ C \cdot O\left(\frac{1}{\sqrt{m}}, \log \frac{1}{\delta}\right) + |a| \cdot O(\|\boldsymbol{\mu}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{S}}\|_2) \end{aligned}$$

Generalization bound

- Loss is a -LOL and Lipschitz
- Bounded feature and model spaces
- Bounded loss $C \doteq \max_x \ell(x)$
- Let $\hat{\boldsymbol{\theta}} \doteq \operatorname{argmin}_{\boldsymbol{\theta} \in \mathcal{H}} R_{\mathcal{S}, \ell}(\boldsymbol{\theta})$

Then for any $\delta > 0$, with probability at least $1 - \delta$:

$$R_{\mathcal{D}, \ell}(\hat{\boldsymbol{\theta}}) - \inf_{\boldsymbol{\theta} \in \mathcal{H}} R_{\mathcal{D}, \ell}(\boldsymbol{\theta}) \leq O\left(\frac{1}{\sqrt{m}}\right) + \text{complexity of space } \mathcal{H}$$

$$C \cdot O\left(\frac{1}{\sqrt{m}}, \log \frac{1}{\delta}\right) + |a| \cdot O(\|\boldsymbol{\mu}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{S}}\|_2)$$

loss non-linearity


$$|a| \cdot O\left(\frac{1}{\sqrt{m}}, \log \frac{1}{\delta}\right)$$

Example: SGD (step 2)

Algorithm SGD

Input: \mathcal{S} , ℓ ;

$$m \leftarrow |\mathcal{S}|$$

$$\boldsymbol{\theta}^0 \leftarrow \mathbf{0}$$

For any $t = 1, 2, \dots$:

 Sample $i \sim U([m])$

$$\eta^t \leftarrow 1/t$$

 Pick any $\mathbf{v} \in \partial\ell(y_i \langle \boldsymbol{\theta}^t, \mathbf{x}_i \rangle)$

$$\boldsymbol{\theta}^{t+1} \leftarrow \boldsymbol{\theta}^t - \eta^t \mathbf{v}$$

Output: $\boldsymbol{\theta}^{t+1}$

Example: SGD (step 2)

Algorithm μ SGD

Input: $\mathcal{S}_{2x}, \mu, \ell$ is a -LOL;

$$m \leftarrow |\mathcal{S}_{2x}|$$

$$\theta^0 \leftarrow 0$$

For any $t = 1, 2, \dots$:

 Sample $i \sim U([m])$

$$\eta^t \leftarrow 1/t$$

 Pick any $v \in \partial\ell(y_i \langle \theta^t, x_i \rangle)$

$$\theta^{t+1} \leftarrow \theta^t - \eta^t(v + a\mu/2)$$

only changes
wrt SGD

Output: θ^{t+1}

Example: SGD (step 2)

Algorithm μ SGD

Input: $\mathcal{S}_{2x}, \mu, \ell$ is a -LOL;

$$m \leftarrow |\mathcal{S}_{2x}|$$

$$\theta^0 \leftarrow 0$$

For any $t = 1, 2, \dots$:

Sample $i \sim U([m])$

$$\eta^t \leftarrow 1/t$$

Pick any $v \in \partial\ell(y_i \langle \theta^t, x_i \rangle)$

$$\theta^{t+1} \leftarrow \theta^t - \eta^t(v + a\mu/2)$$

only changes
wrt SGD

Output: θ^{t+1}

- Similar with proximal algorithms

Outline

- Motivations
- Statistical sufficiency & loss factorization
- **Asymmetric label noise I: theory with linear models**
- Asymmetric label noise II: neural nets
- Other learning settings and conclusion

Asymmetric label noise

Sample $\tilde{\mathcal{S}} = \{(\tilde{x}_i, \tilde{y}_i)\}_{i=1}^m$ corrupted by
asymmetric noise rates p_+, p_-

Asymmetric label noise

Sample $\tilde{\mathcal{S}} = \{(\tilde{x}_i, \tilde{y}_i)\}_{i=1}^m$ corrupted by
asymmetric noise rates p_+, p_-

By the method of [Natarajan et al. '13] an
unbiased estimator of $\mu_{\mathcal{S}}$ is

$$\hat{\mu}_{\mathcal{S}} \doteq \mathbb{E}_{\tilde{\mathcal{S}}} \left[\frac{y - (p_- - p_+)}{1 - p_- - p_+} \mathbf{x} \right], \quad \mathbb{E}_{\tilde{\mathcal{D}}}[\hat{\mu}_{\mathcal{S}}] = \mu_{\mathcal{S}}$$

Asymmetric label noise

Sample $\tilde{\mathcal{S}} = \{(\tilde{x}_i, \tilde{y}_i)\}_{i=1}^m$ corrupted by
asymmetric noise rates p_+, p_-

By the method of [Natarajan et al. '13] an
unbiased estimator of $\mu_{\mathcal{S}}$ is

$$\hat{\mu}_{\mathcal{S}} \doteq \mathbb{E}_{\tilde{\mathcal{S}}} \left[\frac{y - (p_- - p_+)}{1 - p_- - p_+} \mathbf{x} \right], \quad \mathbb{E}_{\tilde{\mathcal{D}}}[\hat{\mu}_{\mathcal{S}}] = \mu_{\mathcal{S}}$$

This is step (1), then run μ -SGD for (2).

Generalization bound under noise

Same as before, except that now we learn with $\hat{\mu}_{\mathcal{S}}$

Then for any $\delta > 0$, with probability at least $1 - \delta$:

$$\begin{aligned} R_{\mathcal{D},\ell}(\hat{\boldsymbol{\theta}}) - \inf_{\boldsymbol{\theta} \in \mathcal{H}} R_{\mathcal{D},\ell}(\boldsymbol{\theta}) &\leq O\left(\frac{1}{\sqrt{m}}\right) + \\ C \cdot O\left(\frac{1}{\sqrt{m}}, \log \frac{1}{\delta}\right) + \frac{|a|}{1 - p_- - p_+} \cdot O\left(\frac{1}{\sqrt{m}}, \log \frac{1}{\delta}\right) \end{aligned}$$

Generalization bound under noise

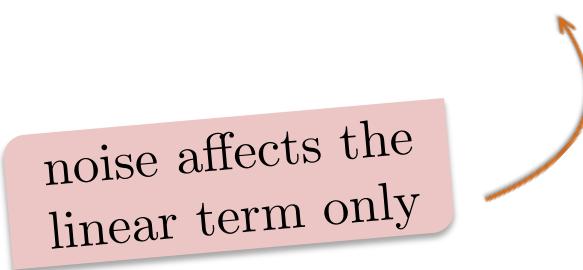
Same as before, except that now we learn with $\hat{\mu}_{\mathcal{S}}$

Then for any $\delta > 0$, with probability at least $1 - \delta$:

$$R_{\mathcal{D},\ell}(\hat{\boldsymbol{\theta}}) - \inf_{\boldsymbol{\theta} \in \mathcal{H}} R_{\mathcal{D},\ell}(\boldsymbol{\theta}) \leq O\left(\frac{1}{\sqrt{m}}\right) + C \cdot O\left(\frac{1}{\sqrt{m}}, \log \frac{1}{\delta}\right) + \frac{|a|}{1 - p_- - p_+} \cdot O\left(\frac{1}{\sqrt{m}}, \log \frac{1}{\delta}\right)$$

complexity
untouched

noise affects the
linear term only



Empirics

- Artificially corrupted data. Noise rates up to ~50%
- SGD vs. μ -SGD with the same parameters
- Test error average difference over 25 runs



$(p_-, p_+) \rightarrow$	(.00, .00)		(.20, .00)		(.20, .10)		(.20, .20)		(.20, .30)		(.20, .40)		(.20, .49)	
dataset	SGD	μ SGD												
australian	0.13	+.01	0.15	-.01	0.14	\pm .00	0.14	+.01	0.16	-.01	0.26	-.09	0.45	-.25
breast-can.	0.02	+.01	0.03	\pm .00	0.03	\pm .00	0.03	\pm .00	0.05	-.01	0.11	-.06	0.17	-.08
diabetes	0.28	-.03	0.29	-.03	0.29	-.03	0.27	-.02	0.28	-.02	0.39	-.13	0.59	-.22
german	0.27	-.02	0.26	\pm .00	0.27	-.02	0.29	-.02	0.31	-.01	0.31	\pm .00	0.31	\pm .00
heart	0.15	+.01	0.17	-.01	0.16	\pm .00	0.17	\pm .00	0.18	-.01	0.26	-.08	0.35	-.15
housing	0.17	-.03	0.23	-.05	0.22	-.04	0.20	-.02	0.22	-.03	0.34	-.12	0.41	-.13
ionosphere	0.14	+.05	0.19	-.05	0.20	-.05	0.20	-.03	0.21	-.03	0.35	-.13	0.54	-.29
sonar	0.27	\pm .00	0.29	+.02	0.29	+.01	0.34	-.04	0.36	-.03	0.43	-.10	0.45	-.05

Empirics

- Artificially corrupted data. Noise rates up to ~50%
- SGD vs. μ -SGD with the same parameters
- Test error average difference over 25 runs

$(p_-, p_+) \rightarrow$	(.00, .00)		(.20, .00)		(.20, .10)		(.20, .20)		(.20, .30)		(.20, .40)		(.20, .49)	
dataset	SGD	μ SGD												
australian	0.13	+.01	0.15	-.01	0.14	\pm .00	0.14	+.01	0.16	-.01	0.26	-.09	0.45	-.25
breast-can.	0.02	+.01	0.03	\pm .00	0.03	\pm .00	0.03	\pm .00	0.05	-.01	0.11	-.06	0.17	-.08
diabetes	0.28	-.03	0.29	-.03	0.29	-.03	0.27	-.02	0.28	-.02	0.39	-.13	0.59	-.22
german	0.27	-.02	0.26	\pm .00	0.27	-.02	0.29	-.02	0.31	-.01	0.31	\pm .00	0.31	\pm .00
heart	0.15	+.01	0.17	-.01	0.16	\pm .00	0.17	\pm .00	0.18	-.01	0.26	-.08	0.35	-.15
housing	0.17	-.03	0.23	-.05	0.22	-.04	0.20	-.02	0.22	-.03	0.34	-.12	0.41	-.13
ionosphere	0.14	+.05	0.19	-.05	0.20	-.05	0.20	-.03	0.21	-.03	0.35	-.13	0.54	-.29
sonar	0.27	\pm .00	0.29	+.02	0.29	+.01	0.34	-.04	0.36	-.03	0.43	-.10	0.45	-.05

=> Still able to learn with one label \sim random

Outline

- Motivations
- Statistical sufficiency & loss factorization
- Asymmetric label noise I: theory with linear models
- **Asymmetric label noise II: neural nets**
- Other learning settings and conclusion

Extended setting

- Multi-class (one-hot): $\mathbf{y} \in \{\mathbf{e}^i : i \in [c]\}$
=> Noise is a transition matrix T
$$T_{ij} = p(\tilde{\mathbf{y}} = \mathbf{e}^j | \mathbf{y} = \mathbf{e}^i)$$
- Estimation of the matrix from noisy data
$$\tilde{\mathcal{S}} \rightarrow \hat{T}$$
- Neural networks

$$\mathbf{h}(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^c$$

Neural network

- For n layers (\mathbf{z} generic input)

$$(\forall i \in [n - 1]) \mathbf{h}^{(i)}(\mathbf{z}) = \sigma(W^{(i)}\mathbf{z} + \mathbf{b}^{(i)}) ,$$

$$\mathbf{h}^{(n)}(\mathbf{z}) = W^{(n)}\mathbf{z} + \mathbf{b}^{(n)}$$

$$\Rightarrow \mathbf{h} = (\mathbf{h}^{(n)} \circ \mathbf{h}^{(n-1)} \circ \dots \circ \mathbf{h}^{(1)})$$

e.g. ReLU,
 $\sigma = \max(0, z)$

- Cross-entropy loss

$$\mathbf{y}^\top \ell(\mathbf{h}(\mathbf{x})) = -\mathbf{y}^\top \log \text{softmax}(\mathbf{h}(\mathbf{x}))$$

No trivial decoupling

- Cross-entropy (linear-odd) factorization

$$\ell(\mathbf{h}(\mathbf{x})) = - (W^n \mathbf{h}^{n-1}(\mathbf{x}) + \mathbf{b}^n) + 1 \cdot \log \sum_{j \in [c]} \exp\{W_j^n \mathbf{h}^{n-1}(\mathbf{x}) + \mathbf{b}_j^n\}$$

label-independent  last layer =
linear model

No trivial decoupling

- Cross-entropy (linear-odd) factorization

$$\ell(\mathbf{h}(\mathbf{x})) = - (W^n \mathbf{h}^{n-1}(\mathbf{x}) + \mathbf{b}^n) + 1 \cdot \log \sum_{j \in [c]} \exp\{W_j^n \mathbf{h}^{n-1}(\mathbf{x}) + \mathbf{b}_j^n\}$$

label-independent  last layer = linear model

- The “sufficient statistic” for the labels is a function of all layers but the last: $\mathbb{E}_{\mathcal{S}}[\mathbf{h}^{n-1}(\mathbf{x})]$

No trivial decoupling

- Cross-entropy (linear-odd) factorization

$$\ell(\mathbf{h}(\mathbf{x})) = - (W^n \mathbf{h}^{n-1}(\mathbf{x}) + \mathbf{b}^n) + 1 \cdot \log \sum_{j \in [c]} \exp\{W_j^n \mathbf{h}^{n-1}(\mathbf{x}) + \mathbf{b}_j^n\}$$

label-independent →

last layer =
linear model

- The “sufficient statistic” for the labels is a function of all layers but the last: $\mathbb{E}_{\mathcal{S}}[\mathbf{h}^{n-1}(\mathbf{x})]$
- But: if we know the noise T , we can make the statistic unbiased **while training**. By [Natarajan et al. '13] :

$$\mathbb{E}_{\mathcal{S}}[T^{-1}(-W^n \mathbf{h}^{n-1}(\mathbf{x}) - \mathbf{b}^n)]$$

is unbiased and so is the whole loss

Noise estimation [Menon et al. '15]

- Train on noisy data and obtain $p(\tilde{y}|\mathbf{x})$
- Then estimate \hat{T} by

$$\forall i, j \quad \begin{cases} \bar{\mathbf{x}}^i = \operatorname{argmax}_{(\mathbf{x}, \cdot) \in \tilde{\mathcal{S}}} p(\tilde{y} = e^i | \mathbf{x}) \\ T_{ij} = p(\tilde{y} = e^j | \bar{\mathbf{x}}^i) \end{cases}$$

- **Hp:** there are some “perfect examples”, and the net models $p(\tilde{y}|\mathbf{x})$ perfectly

Noise estimation [Menon et al. '15]

- Train on noisy data and obtain $p(\tilde{y}|\mathbf{x})$
- Then estimate \hat{T} by

$$\forall i, j \left[\begin{array}{l} \bar{\mathbf{x}}^i = \operatorname{argmax}_{(\mathbf{x}, \cdot) \in \tilde{\mathcal{S}}} p(\tilde{y} = e^i | \mathbf{x}) \\ T_{ij} = p(\tilde{y} = e^j | \bar{\mathbf{x}}^i) \end{array} \right]$$

- **Hp:** there are some “perfect examples”, and the net models $p(\tilde{y}|\mathbf{x})$ perfectly
- **Rational:** mistakes on “perfect examples” must be due to the noise

2-step solution with neural nets

(1) Train the network on $\tilde{\mathcal{S}}$ to obtain \hat{T}

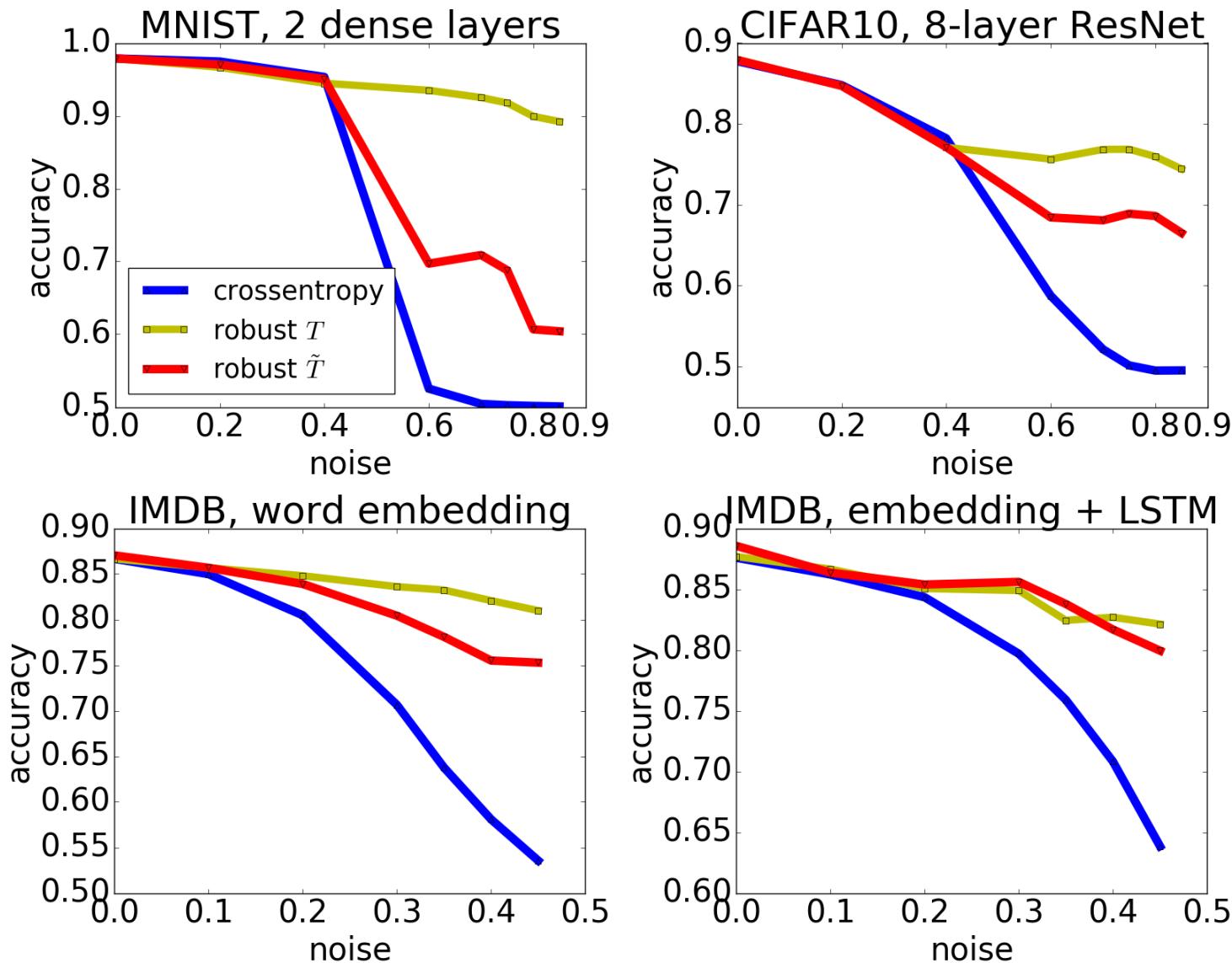
$$\operatorname{argmin}_{\boldsymbol{h}} R_{\tilde{\mathcal{S}}, \ell}(\boldsymbol{h}) \rightarrow p(\tilde{\boldsymbol{y}}|\boldsymbol{x}) \rightarrow \hat{T}$$

(2) Re-train the network correcting the sufficient statistic by \hat{T}^{-1}

$$\boldsymbol{h}^* = \operatorname{argmin}_{\boldsymbol{h}} R_{\tilde{\mathcal{S}}, \ell}(\boldsymbol{h}, \hat{T}^{-1})$$

no change in
back-propagation

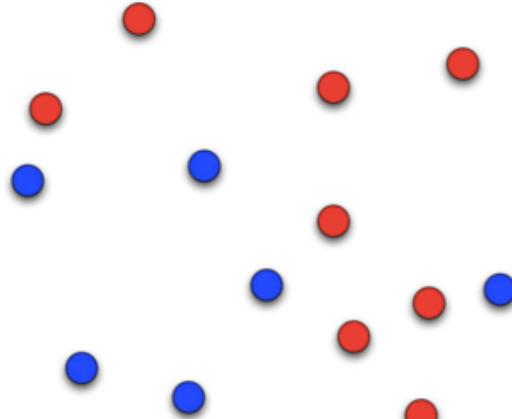
Empirics: inject sparse, asymmetric T



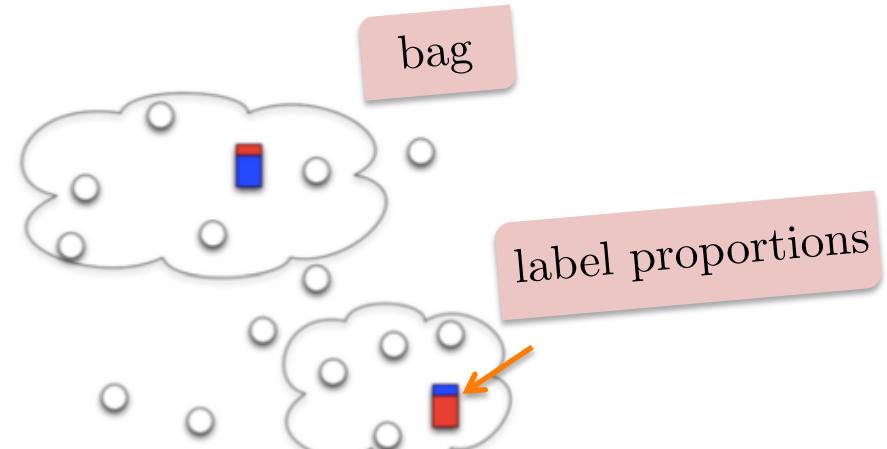
Outline

- Motivations
- Statistical sufficiency & loss factorization
- Asymmetric label noise I: theory with linear models
- Asymmetric label noise II: neural nets
- **Other learning settings and conclusion**

Learning from label proportions

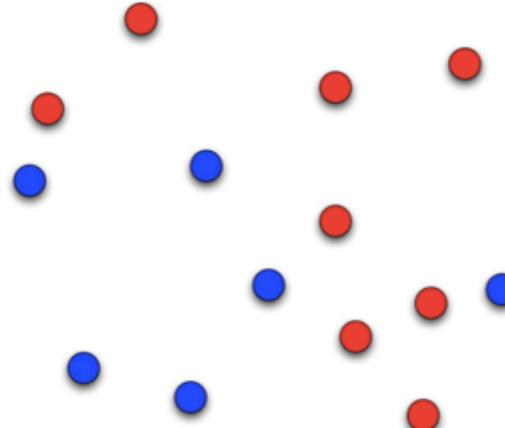


(a) Supervised Learning

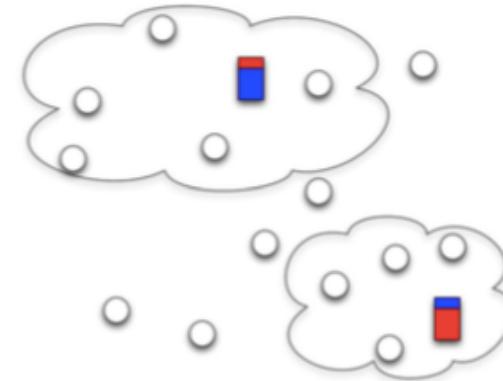


(d) Learning from Proportions

Learning from label proportions



(a) Supervised Learning



(d) Learning from Proportions

Step (1)

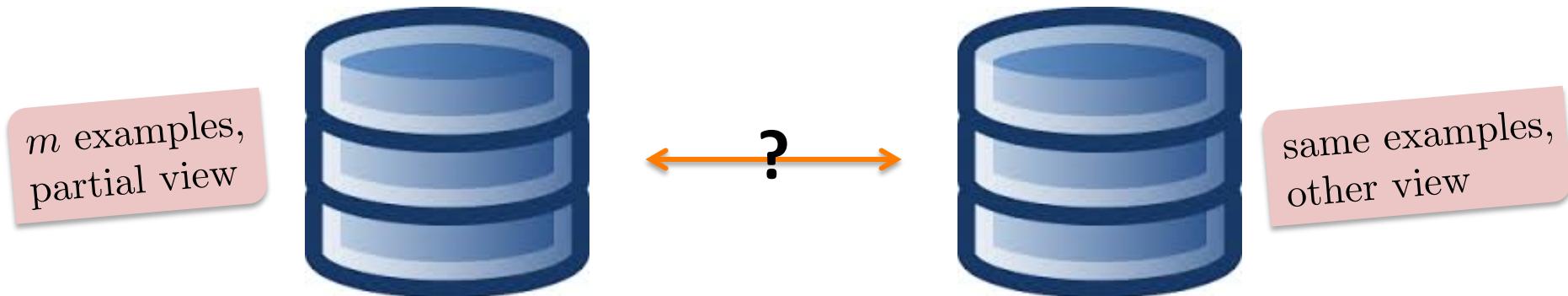
$$\begin{aligned} \hat{U} &= \underset{U}{\operatorname{argmin}} \operatorname{tr} \left((B - \Pi U)^\top (B - \Pi U) \right) + \gamma \operatorname{tr} (U^\top L U) \\ \hat{U} &\rightarrow \hat{\mu} \quad (\text{by linear algebra}) \end{aligned}$$

bag centroids

label proportions

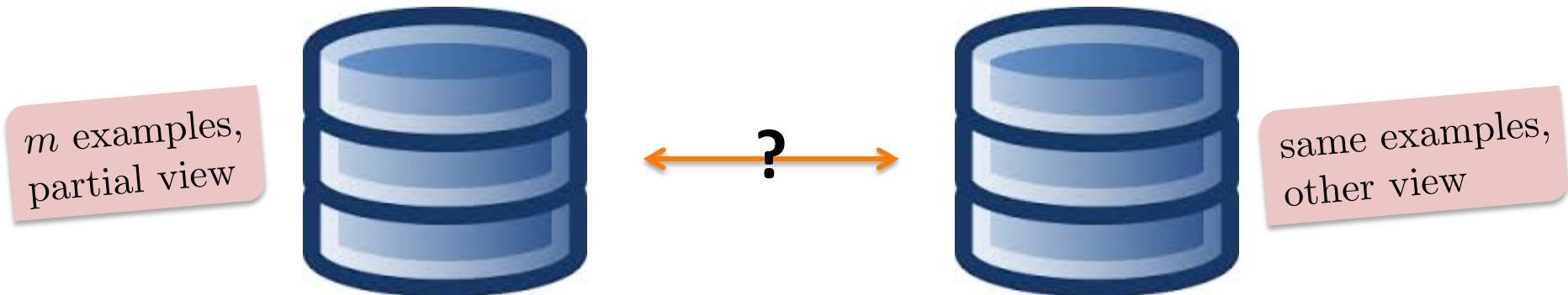
Laplacian regularizer

Learning from distributed datasets



- Data is **vertically partitioned**: same examples, different features
- “who-is-who” not known: **no shared IDs**
- Goal: learn in the product feature space

Learning from distributed datasets



- Data is **vertically partitioned**: same examples, different features
 - “who-is-who” not known: **no shared IDs**
 - Goal: learn in the product feature space
-
- NO: entity matching
 - YES: compute sufficient statistics, Rademacher observations

Conclusion

- Sufficiency is a powerful tool
 - decoupling & modularity
 - abstraction
 - computational saving (by compression)
- Toward ML
 - less like a bag of tricks
 - more like engineering

(see also J. Langford's learning reductions)

References

- G. Patrini, R. Nock, P. Rivera, T. Caetano, **(Almost) no label no cry**, NIPS'14
- G. Patrini, F. Nielsen, R. Nock, M. Carioni, **Loss factorization, weakly supervised learning and label noise robustness**, ICML'16
- G. Patrini, R. Nock, S. Hardy, T. Caetano, **Fast learning from distributed datasets without entity matching**, IJCAI'16
- G. Patrini, A. Rozza, R. Nock, A. Menon, L. Qu, **Making Neural Networks Robust to Label Noise: a Loss Correction Approach**, NIPS'16 (submitted)

Co-authored

- R. Nock, G. Patrini, A. Friedman, **Rademacher observations, private data and boosting**, ICML'15

Patent applications

- R. Nock, G. Patrini, T. Caetano, **Learning with transformed data**, 2015
- R. Nock, G. Patrini, **Learning with distributed data**, 2016