

# On Structured Prediction Theory with Calibrated Convex Surrogate Losses

Anton Osokin



Francis Bach



Simon Lacoste-Julien



08.09.2017

# Structured prediction

- Data is often very structured
  - sequences, tables, images, video, etc.
- Structured prediction = output is structured

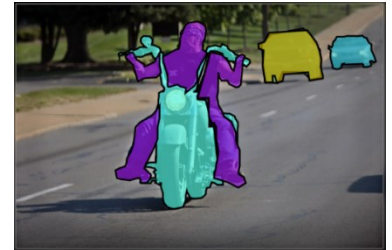
- Examples:

- Image segmentation

Input data



Output labels



- Handwriting recognition



command

# Inference and learning

- Inference (prediction)

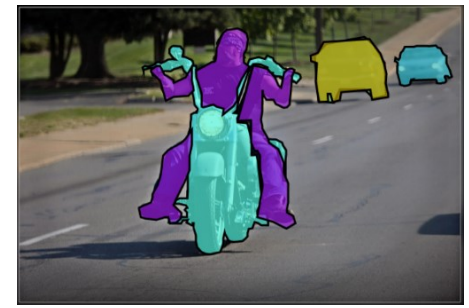


Input data

+

Model

=



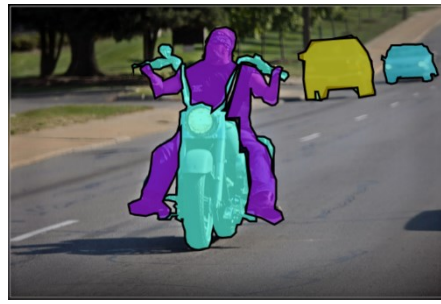
Output labels

- Learning



Input data

+



Output labels

=

Model

Training data

We want a method that works!  
and has guarantees!

Number of iterations to get  $\varepsilon$ -accuracy on the test set

# Conditional likelihood (e.g., CRFs)

Input  $\mathbf{x}$ , output  $\mathbf{y}$ , loss  $L$ ; set of possible  $\mathbf{y}$  is exponential!

Training – maximize conditional likelihood

$$\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} \log p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta})$$

Prediction – minimize expected loss

$$\operatorname{pred}(\mathbf{x}) = \operatorname{argmin}_{\hat{\mathbf{y}}} \sum_{\mathbf{y}} L(\hat{\mathbf{y}}, \mathbf{y}) p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta})$$

Problems:

- hard to connect to prediction error
- training does not know about the loss
- the whole distribution is hard to get

# Empirical risk minimization

Input  $\mathbf{x}$ , output  $\mathbf{y}$ , loss  $L$ ; set of possible  $\mathbf{y}$  is exponential!  
We want to minimize the population risk

$$\mathcal{R}_L(\mathbf{f}) := \mathbf{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} L\left(\operatorname{argmax}(\mathbf{f}(\mathbf{x})), \mathbf{y}\right)$$

Here  $\mathbf{f}$  gives a score for each label

Usually we cannot optimize the risk, so we use a surrogate

$$\mathcal{R}_\Phi(\mathbf{f}) := \mathbf{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} \Phi(\mathbf{y}, \mathbf{f}(\mathbf{x}))$$

Examples:

$$\Phi_{\text{SSVM}}(\mathbf{f}, \mathbf{y}) := \max_{\hat{\mathbf{y}} \in \mathcal{Y}} (f_{\hat{\mathbf{y}}} + L(\hat{\mathbf{y}}, \mathbf{y})) - f_{\mathbf{y}}$$

$$\Phi_{\log}(\mathbf{f}, \mathbf{y}) := -f_{\mathbf{y}} + \log \sum_{\hat{\mathbf{y}} \in \mathcal{Y}} \exp f_{\hat{\mathbf{y}}}$$

# Where is structure?

Input  $\mathbf{x}$ , output  $\mathbf{y}$ , loss  $L$ ; set of possible  $\mathbf{y}$  is exponential!  
We want to minimize the population risk

$$\mathcal{R}_L(\mathbf{f}) := \mathbf{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} L\left(\operatorname{argmax}(\mathbf{f}(\mathbf{x})), \mathbf{y}\right)$$

Here  $\mathbf{f}$  gives a score for each label

Usually we cannot optimize the risk, so we use a surrogate

$$\mathcal{R}_\Phi(\mathbf{f}) := \mathbf{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} \Phi(\mathbf{y}, \mathbf{f}(\mathbf{x}))$$

- 1) Loss  $L$  is structured (Hamming, block 0-1)
- 2) Scores are structured:  $\mathbf{f} \in \mathcal{F}$  (for example, separable)

# Empirical risk minimization

Input  $\mathbf{x}$ , output  $\mathbf{y}$ , loss  $L$ ; set of possible  $\mathbf{y}$  is exponential!  
We want to minimize the population risk

$$\mathcal{R}_L(\mathbf{f}) := \mathbf{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} L(\operatorname{argmax}(\mathbf{f}(\mathbf{x})), \mathbf{y})$$

Here  $\mathbf{f}$  gives a score for each label

Usually we cannot optimize the risk, so we use a surrogate

$$\mathcal{R}_\Phi(\mathbf{f}) := \mathbf{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} \Phi(\mathbf{y}, \mathbf{f}(\mathbf{x}))$$

Issues to be dealt with:

- Surrogate  $\Phi$  instead of the actual loss  $L$
- Finite dataset
- Optimization accuracy

Consistency

} Online SGD



# Consistency for binary classification

Input  $\mathbf{x}$ , output  $y \in \{-1, 1\}$ , loss  $L(\hat{y}, y) = [\hat{y} \neq y]$

We want to minimize the population risk

$$\mathcal{R}_L(f) := \mathbf{E}_{(x,y) \sim \mathcal{D}} [yf(x) > 0]$$

Surrogate risk:  $\mathcal{R}_\Phi(f) := \mathbf{E}_{(x,y) \sim \mathcal{D}} \Phi(yf(x))$

**Theorem** (Bartlett et al., 2006)

Assume  $\Phi$  is convex, differentiable at 0 and  $\Phi'(0) < 0$ .

Then  $\Phi$  is classification-calibrated (consistent).

Examples: exponential, hinge, logistic, truncated quadratic

# Calibration function

Calibration function connects the surrogate and the loss.

$$\begin{aligned} H_{\Phi, L, \mathcal{F}}(\varepsilon) &= \min_{\mathbf{f}, \mathbf{q}} \delta\phi(\mathbf{q}, \mathbf{f}) \\ \text{s.t. } &\delta\ell(\mathbf{q}, \mathbf{f}) \geq \varepsilon, \\ &\mathbf{f} \in \mathcal{F}, \mathbf{q} \in \Delta_k, \end{aligned}$$

With expected and excess loss/surrogate defined as

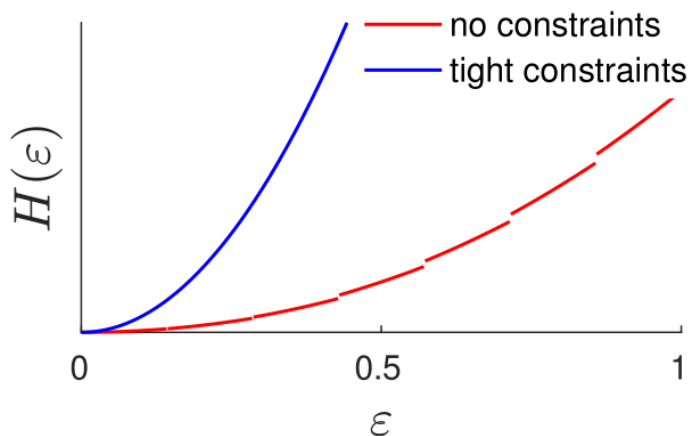
$$\begin{aligned} \ell(\mathbf{q}, \mathbf{f}) &= \sum_{c=1}^k q_c L(\operatorname{argmax}(\mathbf{f}), c), \quad \phi(\mathbf{q}, \mathbf{f}) = \sum_{c=1}^k q_c \Phi(c, \mathbf{f}) \\ \delta\phi(\mathbf{q}, \mathbf{f}) &= \phi(\mathbf{q}, \mathbf{f}) - \min_{\hat{\mathbf{f}} \in \mathcal{F}} \phi(\mathbf{q}, \hat{\mathbf{f}}) \\ \delta\ell(\mathbf{q}, \mathbf{f}) &= \ell(\mathbf{q}, \mathbf{f}) - \min_{\hat{\mathbf{f}} \in \mathcal{F}} \ell(\mathbf{q}, \hat{\mathbf{f}}) \end{aligned}$$

(Zhang, 2004) and others

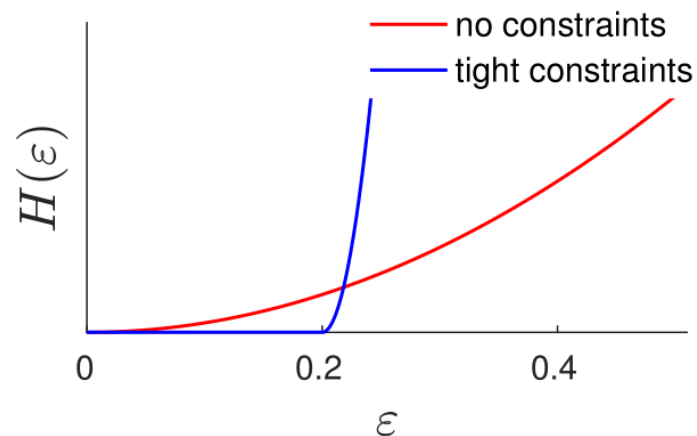
# Calibration function

Calibration function connects the surrogate and the loss.

$$\begin{aligned} H_{\Phi, L, \mathcal{F}}(\varepsilon) &= \min_{\mathbf{f}, \mathbf{q}} \delta\phi(\mathbf{q}, \mathbf{f}) \\ \text{s.t. } &\delta\ell(\mathbf{q}, \mathbf{f}) \geq \varepsilon, \\ &\mathbf{f} \in \mathcal{F}, \mathbf{q} \in \Delta_k, \end{aligned}$$



(a): Hamming loss



(b): Mixed loss

(Zhang, 2004) and others

# Calibration function and consistency

Calibration function connects the surrogate and the loss.

$$\begin{aligned} H_{\Phi, L, \mathcal{F}}(\varepsilon) &= \min_{\mathbf{f}, \mathbf{q}} \delta\phi(\mathbf{q}, \mathbf{f}) \\ \text{s.t. } &\delta\ell(\mathbf{q}, \mathbf{f}) \geq \varepsilon, \\ &\mathbf{f} \in \mathcal{F}, \mathbf{q} \in \Delta_k, \end{aligned}$$

If we optimize up to  $\mathcal{R}_{\Phi}(\mathbf{f}) < \mathcal{R}_{\Phi}^* + H_{\Phi, L, \mathcal{F}}(\varepsilon)$

We get  $\mathcal{R}_L(\mathbf{f}) < \mathcal{R}_L^* + \varepsilon$

If  $H_{\Phi, L, \mathcal{F}}(\varepsilon) > 0$ ,  $\varepsilon > 0$  the surrogate is consistent.

(Zhang, 2004) and others

# If consistent, what can go wrong?

Calibration function connects the surrogate and the loss.

$$\begin{aligned} H_{\Phi, L, \mathcal{F}}(\varepsilon) &= \min_{\mathbf{f}, \mathbf{q}} \delta\phi(\mathbf{q}, \mathbf{f}) \\ \text{s.t. } &\delta\ell(\mathbf{q}, \mathbf{f}) \geq \varepsilon, \\ &\mathbf{f} \in \mathcal{F}, \mathbf{q} \in \Delta_k, \end{aligned}$$

If we optimize up to  $\mathcal{R}_{\Phi}(\mathbf{f}) < \mathcal{R}_{\Phi}^* + H_{\Phi, L, \mathcal{F}}(\varepsilon)$

We get  $\mathcal{R}_L(\mathbf{f}) < \mathcal{R}_L^* + \varepsilon$

If  $H_{\Phi, L, \mathcal{F}}(\varepsilon) > 0$ ,  $\varepsilon > 0$  the surrogate is consistent.

# If consistent, what can go wrong?

Calibration function connects the surrogate and the loss.

$$\begin{aligned} H_{\Phi, L, \mathcal{F}}(\varepsilon) &= \min_{\mathbf{f}, \mathbf{q}} \delta\phi(\mathbf{q}, \mathbf{f}) \\ \text{s.t. } &\delta\ell(\mathbf{q}, \mathbf{f}) \geq \varepsilon, \\ &\mathbf{f} \in \mathcal{F}, \mathbf{q} \in \Delta_k, \end{aligned}$$

Required accuracy  $H_{\Phi, L, \mathcal{F}}(\varepsilon)$  can be exponentially small!

We will never reach it: finite dataset and runtime

# Easy-to-analyze surrogate

Calibration function connects the surrogate and the loss.

$$\begin{aligned} H_{\Phi, L, \mathcal{F}}(\varepsilon) &= \min_{\mathbf{f}, \mathbf{q}} \delta\phi(\mathbf{q}, \mathbf{f}) \\ \text{s.t. } &\delta\ell(\mathbf{q}, \mathbf{f}) \geq \varepsilon, \\ &\mathbf{f} \in \mathcal{F}, \mathbf{q} \in \Delta_k, \end{aligned}$$

Simple surrogate, consistent for any given loss:

$$\Phi_{\text{quad}}(\mathbf{y}, \mathbf{f}) = \frac{1}{2k} \|\mathbf{f} + L(:, \mathbf{y})\|_2^2 \quad \text{where} \quad \mathbf{f}(\mathbf{y}) = F\boldsymbol{\theta}(\mathbf{x})$$

# 01-loss

Calibration function connects the surrogate and the loss.

$$\begin{aligned} H_{\Phi, L, \mathcal{F}}(\varepsilon) &= \min_{\mathbf{f}, \mathbf{q}} \delta\phi(\mathbf{q}, \mathbf{f}) \\ \text{s.t. } &\delta\ell(\mathbf{q}, \mathbf{f}) \geq \varepsilon, \\ &\mathbf{f} \in \mathcal{F}, \mathbf{q} \in \Delta_k, \end{aligned}$$

$k$  – number of classes

Calibration function:  $H_{\Phi_{\text{quad}}, L_{01}, \mathbb{R}^k}(\varepsilon) = \frac{\varepsilon^2}{4k}$

If  $k$  is big then  $H$  is very small, i.e. we need crazy accuracy!



# Block 01-loss

Calibration function connects the surrogate and the loss.

$$\begin{aligned} H_{\Phi, L, \mathcal{F}}(\varepsilon) &= \min_{\mathbf{f}, \mathbf{q}} \delta\phi(\mathbf{q}, \mathbf{f}) \\ \text{s.t. } &\delta\ell(\mathbf{q}, \mathbf{f}) \geq \varepsilon, \\ &\mathbf{f} \in \mathcal{F}, \mathbf{q} \in \Delta_k, \end{aligned}$$

$k$  – number of classes,  $s$  – block size

Calibration function:  $H_{\Phi_{\text{quad}}, L, \mathbb{R}^k}(\varepsilon) = \frac{\varepsilon^2}{4k} \frac{2s}{s+1} \leq 2 \frac{\varepsilon^2}{4k}$

The loss is highly structured, but  $H$  is still very small!

What is missing?

# Block 01-loss + constraints on the scores

Calibration function connects the surrogate and the loss.

$$\begin{aligned} H_{\Phi, L, \mathcal{F}}(\varepsilon) &= \min_{\mathbf{f}, \mathbf{q}} \delta\phi(\mathbf{q}, \mathbf{f}) \\ \text{s.t. } \delta\ell(\mathbf{q}, \mathbf{f}) &\geq \varepsilon, \\ \mathbf{f} &\in \mathcal{F}, \mathbf{q} \in \Delta_k, \end{aligned}$$

$k$  – number of classes,  $s$  – block size,  $d$  – number of blocks

Calibration function:  $H_{\Phi_{\text{quad}}, L, \mathcal{L}}(\varepsilon) = \frac{\varepsilon^2}{4d}$   **$d$  is small!**

$d$  represents complexity of the loss

# Hamming loss

Hamming loss on binary sequences:

$$H_{\Phi_{\text{quad}}, \text{Ham}, \mathbb{R}^k} = \frac{\varepsilon^2}{4k} \text{ for small } \varepsilon \quad \text{very small!}$$

With constraints on the scores:

$$H_{\Phi_{\text{quad}}, \text{Ham}, \mathcal{F}} = \frac{\varepsilon^2}{8T} \quad \text{much better!}$$

$T$  – length of the sequence

# Lower bound for any loss

**Theorem 1.** *For any loss  $L$ , its quadratic surrogate  $\Phi_{\theta}$ , and a score subspace  $\mathcal{F}$  containing the column space of  $L$ , the calibration function can be lower bounded:*

$$H_{\Phi_{\theta}, L, \mathcal{F}}(\varepsilon) \geq \frac{\varepsilon^2}{2k} \min_{i,j=1,\dots,k} \frac{1}{\|P_{\mathcal{F}}\delta_{ij}\|_2^2} \geq \frac{\varepsilon^2}{4k}$$

where  $P_{\mathcal{F}}$  is the projection of subspace  $\mathcal{F}$  of allowed scores,  $\delta_{ij} = \delta_i - \delta_j \in \mathbb{R}^k$  defines the error direction.

The bound is tight in some cases!

(Block 01 loss, Hamming loss with appropriate constraints)

The bound suggests that good losses are low-rank.

Subspace should not be aligned with  $\delta_{ij}$

OK, say calibration function is large.  
Are we done?

No, scale of calibration function is defined arbitrarily.

We have to connect to optimization and trace constants.

# Minimizing the surrogate risk

Consider minimizing

$$\mathcal{R}_{\Phi}(\mathbf{f}) = \mathbf{E}_{\mathbf{x}, \mathbf{y} \sim \mathcal{D}} \frac{1}{2k} \|\mathbf{f}(\mathbf{x}) + L(:, \mathbf{y})\|_2^2$$

Where we assume linear dependency on features  $\psi(\mathbf{x}) \in \mathbb{R}^d$

$$\mathbf{f}(\mathbf{x}) = FW\psi(\mathbf{x}), \quad W \in \mathbb{R}^{r \times d}$$

Stochastic gradient:

$$\nabla_W \mathcal{R}_{\Phi}^{(n)} = \frac{1}{k} F^{\top} (FW\psi(\mathbf{x}) - L(:, \mathbf{y})) \psi(\mathbf{x})^{\top}$$

Now we just need a suitable algorithm with convergence rate.

# Normalizing by SGD

Projected averaged SGD with constant step size

$$W^{(n)} := P_D [W^{(n-1)} - \gamma \nabla_W \mathcal{R}_\Phi^{(n)}], \quad \gamma = \frac{2D}{M\sqrt{N}}$$

$D$  – bounds the distance to optimum

$M^2$  – bounds the expectation of squared L2-norm of the gradient

(Nemirovski et al.) Convergence rate of averaged iterates

$$\mathbf{E}[\mathcal{R}_\Phi(\bar{\mathbf{f}}^{(N)})] - \mathcal{R}_\Phi^* \leq \frac{2DM}{\sqrt{N}}, \quad \bar{\mathbf{f}}^{(N)} := \frac{1}{N} \sum_{n=1}^N F W^{(n)} \psi(\mathbf{x})^\top$$

Combining this with calibration functions?

# Big assumption: well-specified optimization

There exists a global minimum of  $\mathcal{R}_\Phi(f)$   
w.r.t. all measurable function  $f(\mathbf{x})$   
that belongs to the function class  $f(\mathbf{x}) = FW\psi(\mathbf{x})$

This assumption can be relaxed if we use a universal kernel  
(that can approximate any function)  
and if we are very careful with optimization



# Normalizing by SGD

Projected averaged SGD with constant step size

$$W^{(n)} := P_D [W^{(n-1)} - \gamma \nabla_W \mathcal{R}_\Phi^{(n)}], \quad \gamma = \frac{2D}{M\sqrt{N}}$$

$D$  – bounds the distance to optimum

$M^2$  – bounds the expectation of squared L2-norm of the gradient

(Nemirovski et al.) Convergence rate of averaged iterates

$$\mathbf{E}[\mathcal{R}_\Phi(\bar{\mathbf{f}}^{(N)})] - \mathcal{R}_\Phi^* \leq \frac{2DM}{\sqrt{N}}, \quad \bar{\mathbf{f}}^{(N)} := \frac{1}{N} \sum_{n=1}^N F W^{(n)} \psi(\mathbf{x})^\top$$

Combining this with calibration functions, we need

$$T \geq \frac{4D^2 M^2}{H_{\Phi, L, \mathcal{F}}^2(\varepsilon)} \quad \text{iterations to reach} \quad \mathcal{R}_L(\tilde{\mathbf{f}}^{(T)}) < \mathcal{R}_L^* + \varepsilon,$$

# Normalizing by SGD: quadratic

We need  $T \geq \frac{4D^2 M^2}{H_{\Phi, L, \mathcal{F}}^2(\varepsilon)}$  iterations to reach  $\mathcal{R}_L(\tilde{\mathbf{f}}^{(T)}) < \mathcal{R}_L^* + \varepsilon$

$D$  – bounds the distance to optimum

$M^2$  – bounds the expectation of squared L2-norm of the gradient

For quadratic surrogate

$$DM = L_{\max}^2 \xi(\kappa(F) \sqrt{r} R Q_{\max}), \quad \xi(z) = z^2 + z,$$

$\kappa(F)$  – the condition number of  $F$

$r$  – rank of  $F$

$R$  – upper bound on  $\|\psi(\mathbf{x})\|_2$

$Q_{\max}$  – upper bound on the sum of the norms of marginal probs

# Normalizing by SGD: quadratic

We need  $T \geq \frac{4D^2 M^2}{H_{\Phi, L, \mathcal{F}}^2(\varepsilon)}$  iterations to reach  $\mathcal{R}_L(\tilde{\mathbf{f}}^{(T)}) < \mathcal{R}_L^* + \varepsilon$

$D$  – bounds the distance to optimum

$M^2$  – bounds the expectation of squared L2-norm of the gradient

For quadratic surrogate

$$DM = L_{\max}^2 \xi(\kappa(F) \sqrt{r} R Q_{\max}), \quad \xi(z) = z^2 + z,$$

**0-1 loss:**  $DM = O(k)$

**Block 0-1 loss:**  $DM = O(b)$

**Hamming loss:**  $DM = O(\log_2^3 k)$

# Conclusion

- What are good surrogates?
- Consistency is not enough
- Calibration functions can be computed
- Scale can be chosen by connection to optimization
- Need to be careful with constants

Future work:

- Bounding calibration functions and constants for more losses
- Generalize analysis to other surrogates
- Make this practical!