

Generic Structure learning problems

we have something we solve by

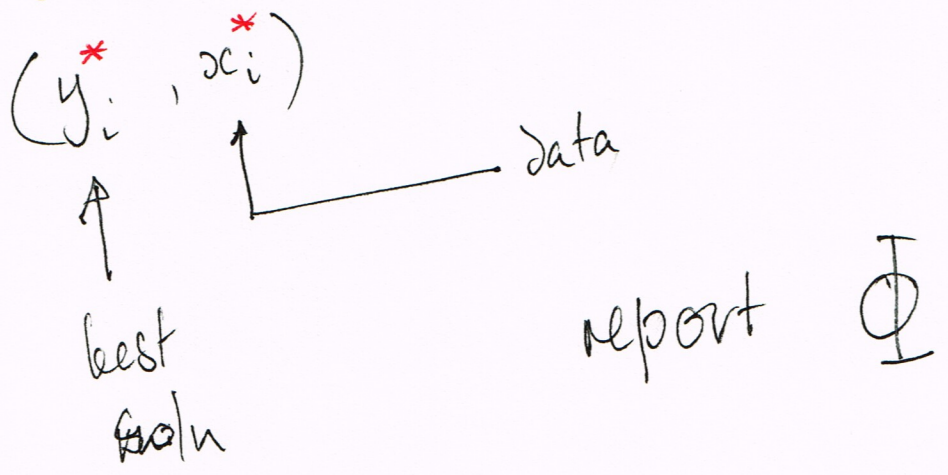
$$\hat{y} = \operatorname{argmax}_y \Phi(y, x)$$

↑
Data.

Many examples:

- matching
- inference for HMMs
- etc.

Now given



This must be hard; but framework

(15)
(Σ)

assume $\Phi(y, x)$

$$= \sum_i \omega_i \phi_i(y, x)$$

Q: how do we get ω ?
↑ "features"

have data so

$$\omega^T \phi(y_i^*, x_i^*) \geq \omega^T \phi(y, x_i^*)$$

This is NOT JUST 1 Ineq.!
TRUE for all y !

- observed y_i^* is better than all others.
- Typically, this isn't enough.

if y "far" from y_i^*

want $\omega^T \phi(y_i^*, x_i^*)$ a lot
bigger than $\omega^T \phi(y, x_i^*)$.

so

$$\omega^T \phi(y_i^*, x_i^*) \geq \omega^T \phi(y, x_i^*) + \frac{\lambda}{2} \|y - y_i^*\|^2$$

↑
or some norm.



true for any y .

This constraint might be violated.

consider

$$\xi_i(y) = \max \left[\omega^T \phi(y, x_i^*) + \frac{\lambda}{2} \|y - y_i^*\|^2 - \omega^T \phi(y_i^*, x_i^*), 0 \right]$$

↑
= 0 if constraint satisfied
> 0 otherwise.

We are interested in worst constraint violation

so

$$\xi_i = \max_y \left[\max \left(\omega^T \phi(y, x_i^*) + \frac{\lambda}{2} \|y - y_i^*\|^2 - \omega^T \phi(y_i^*, x_i^*), 0 \right) \right]$$

Now imagine we have many examples

17
4

- small ω is a good idea
- So

$$\min_{\omega} \sum_i \xi_i + \frac{\mu}{2} \omega^T \omega$$

$$\text{st } \xi_i = \max_y \left[\max \left(\omega^T \phi(y, x_i^*) + \frac{\lambda}{2} \|y - y_i^*\|^2 - \omega^T \phi(y_i^*, x_i) \right), 0 \right]$$

Generically, called structure learning

Startlingly, it's convex

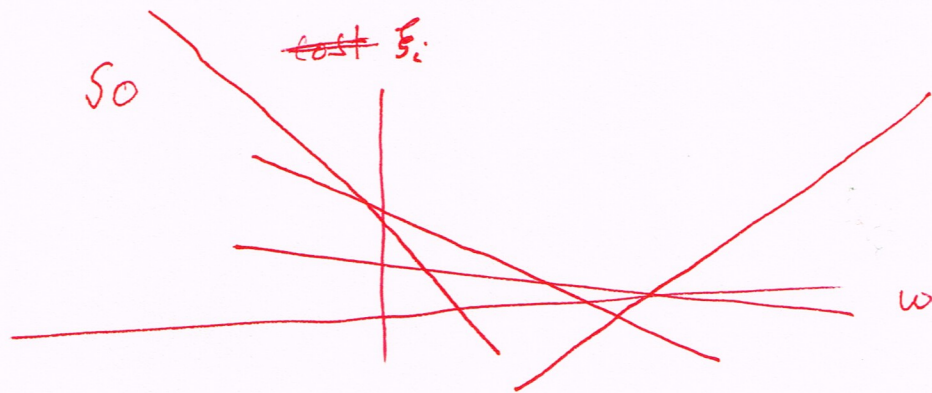
- envelope of convex fns.

- Notice

- for any particular example

- if the max doesn't change, it's linear

- and we're minimizing



- Structure learning is NASTY

- Variety of strategies

- subgradient alg.

$w^0 = \text{init}$

- • obtain subgradient ∂cost
- $w^{(c+1)} = w^{(c)} - \eta \partial$

- This will converge but

- subgradient is hugely expensive
 - must do inference at every example

- Stochastic subgradient

- • obtain stoch est of subgradient (draw sample)
- do descent

- Notice if there were few examples, you could get clever

eg. choose samp using ξ_i etc