

Structure and Support Vector Machines

Outline

- SVMs for structured outputs
 - Declarative view
 - Procedural view

Notation for Linear Models

- Training data: $\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$
- Testing data: $\{(x_{N+1}, y_{N+1}), \dots (x_{N+N'}, y_{N+N'})\}$
- Feature function: \mathbf{g}
- Weights: \mathbf{w}
- Decoding:

$$\text{decode}(\mathbf{w}, \mathbf{x}) = \arg \max_{\mathbf{y}} \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y})$$

- Learning:

$$\text{learn}(\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N) = \arg \max_{\mathbf{w}} \Phi(\mathbf{w}, \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N)$$

- Evaluation:

$$\frac{1}{N'} \sum_{i=1}^{N'} \text{cost}(\text{decode}(\text{learn}(\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N), \mathbf{x}_{N+i}), \mathbf{y}_{N+i})$$

Empirical Risk Minimization

- A unifying framework for many learning algorithms.

$$\begin{aligned}\text{learn} \left(\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N \right) &= \arg \max_{\mathbf{w}} \Phi \left(\mathbf{w}, \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N \right) \\ &= \arg \min_{\mathbf{w}} \underbrace{\frac{1}{N} \sum_{i=1}^N L(\mathbf{w}, \mathbf{x}_i, \mathbf{y}_i) + R(\mathbf{w})}_{\approx \mathbb{E}[L(\mathbf{w}, \mathbf{X}, \mathbf{Y})]}\end{aligned}$$

- Many options for the loss function L and the regularization function R .

Loss Functions You Know

Name	Expression of $L(\mathbf{w}, \mathbf{x}, \mathbf{y})$
Log loss (joint)	$-\log p(\mathbf{x}, \mathbf{y} \mid \mathbf{w})$
Log loss (conditional)	$-\log p(\mathbf{y} \mid \mathbf{x}, \mathbf{w})$
Zero-one loss	$\mathbf{1}\{\text{decode}(\mathbf{w}, \mathbf{x}) \neq \mathbf{y}\}$
Expected zero-one loss	$1 - p(\mathbf{y} \mid \mathbf{x}, \mathbf{w})$

Loss Functions You Know

Name	Expression of $L(\mathbf{w}, \mathbf{x}, \mathbf{y})$
Log loss (joint)	$-\log p(\mathbf{x}, \mathbf{y} \mid \mathbf{w})$
Log loss (conditional)	$-\log p(\mathbf{y} \mid \mathbf{x}, \mathbf{w})$
Cost	$\text{cost}(\text{decode}(\mathbf{w}, \mathbf{x}), \mathbf{y})$
Expected cost, a.k.a. “risk”	$\mathbb{E}_{p(\mathbf{Y} \mid \mathbf{x}, \mathbf{w})} [\text{cost}(\mathbf{Y}, \mathbf{y})]$

Structured Perceptron

- Described as an online algorithm.
- On each iteration, take one example, and update the weights according to:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha (\mathbf{g}(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{g}(\mathbf{x}_t, \text{decode}(\mathbf{w}, \mathbf{x}_t)))$$

Loss Functions You Know

Name	Expression of $L(\mathbf{w}, x, y)$
Log loss (joint)	$-\log p(\mathbf{x}, \mathbf{y} \mid \mathbf{w})$
Log loss (conditional)	$-\log p(\mathbf{y} \mid \mathbf{x}, \mathbf{w})$
Cost	$\text{cost}(\text{decode}(\mathbf{w}, \mathbf{x}), \mathbf{y})$
Expected cost, a.k.a. “risk”	$\mathbb{E}_{p(\mathbf{Y} \mid \mathbf{x}, \mathbf{w})} [\text{cost}(\mathbf{Y}, \mathbf{y})]$
Perceptron loss	$-\mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}) + \max_{\mathbf{y}'} \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}')$

The Ideal Loss Function

- Convex
- Continuous
- Cost-aware

Cost and Margin

- The “margin” is an important concept when we take the linear models point of view.
 - A “large margin” means that the correct output is well-separated from the incorrect outputs.
- Neither log loss nor “perceptron loss” takes into account the *cost* function, though.
 - In other words, some incorrect outputs are worse than others.

Multiclass SVM (Crammer and Singer, 2001)

$$\begin{aligned} \max_{\mathbf{w}} \quad & \gamma \\ \text{s.t.} \quad & \|\mathbf{w}\| \leq 1 \\ & \forall i, \forall \mathbf{y}, \mathbf{w}^\top \mathbf{g}(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{w}^\top \mathbf{g}(\mathbf{x}_i, \mathbf{y}) \geq \begin{cases} \gamma & \text{if } \mathbf{y} \neq \mathbf{y}_i \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

- The above can be understood as a 0-1 cost; let's generalize a bit:

$$\begin{aligned} \max_{\mathbf{w}} \quad & \gamma \\ \text{s.t.} \quad & \|\mathbf{w}\| \leq 1 \\ & \forall i, \forall \mathbf{y}, \mathbf{w}^\top \mathbf{g}(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{w}^\top \mathbf{g}(\mathbf{x}_i, \mathbf{y}) \geq \gamma \text{cost}(\mathbf{y}, \mathbf{y}_i) \end{aligned}$$

Max-Margin Markov Networks

- Starting point: multiclass SVM (Crammer and Singer, 2001)

$$\begin{aligned} \max_{\mathbf{w}} \quad & \gamma \\ \text{s.t.} \quad & \|\mathbf{w}\| \leq 1 \\ & \forall i, \forall \mathbf{y}, \mathbf{w}^\top \mathbf{g}(x_i, \mathbf{y}_i) - \mathbf{w}^\top \mathbf{g}(x_i, \mathbf{y}) \geq \gamma \text{cost}(\mathbf{y}, \mathbf{y}_i) \end{aligned}$$

Max-Margin Markov Networks

- Standard transformation to get rid of explicit mention of γ , plus slack variables in case the constraints cannot be met:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & \forall i, \forall \mathbf{y}, \quad \mathbf{w}^\top \mathbf{g}(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{w}^\top \mathbf{g}(\mathbf{x}_i, \mathbf{y}) \geq \text{cost}(\mathbf{y}, \mathbf{y}_i) - \xi_i \end{aligned}$$

- Notice:

$$\begin{aligned} \forall i, \forall \mathbf{y}, \quad \xi_i &\geq -\mathbf{w}^\top \mathbf{g}(\mathbf{x}_i, \mathbf{y}_i) + \mathbf{w}^\top \mathbf{g}(\mathbf{x}_i, \mathbf{y}) + \text{cost}(\mathbf{y}, \mathbf{y}_i) \\ \forall i, \quad \xi_i &\geq \max_{\mathbf{y}} -\mathbf{w}^\top \mathbf{g}(\mathbf{x}_i, \mathbf{y}_i) + \mathbf{w}^\top \mathbf{g}(\mathbf{x}_i, \mathbf{y}) + \text{cost}(\mathbf{y}, \mathbf{y}_i) \end{aligned}$$

Max-Margin Markov Networks

- Having solved for the slack variables, we can plug in; we now have an unconstrained problem:

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^N -\mathbf{w}^\top \mathbf{g}(\mathbf{x}_i, \mathbf{y}_i) + \max_{\mathbf{y}} \mathbf{w}^\top \mathbf{g}(\mathbf{x}_i, \mathbf{y}) + \text{cost}(\mathbf{y}, \mathbf{y}_i)$$

- Ratliff, Bagnell, and Zinkevich (2007): subgradient descent (or stochastic version) – much, much simpler approach to optimizing this function.
 - And more perceptron-like!

$$-g_j(\mathbf{x}, \mathbf{y}) + g_j(\mathbf{x}, \text{cost_augmented_decode}(\mathbf{w}, \mathbf{x}))$$

Structured Hinge Loss

- Small change to the perceptron loss:

$$L(\mathbf{w}, \mathbf{x}, \mathbf{y}) = -\mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}) + \max_{\mathbf{y}'} \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}') + \text{cost}(\mathbf{y}', \mathbf{y})$$

- Resulting subgradient:

$$-g_j(\mathbf{x}, \mathbf{y}) + g_j(\mathbf{x}, \text{cost_augmented_decode}(\mathbf{w}, \mathbf{x}))$$

- Rather than merely decoding, find a candidate \mathbf{y}' that is both high-scoring and *dangerous*.

Cost-Augmented Decoding

$$\text{decode}(\mathbf{w}, \mathbf{x}) = \arg \max_{\mathbf{y}'} \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}')$$

$$\text{cost_augmented_decode}(\mathbf{w}, \mathbf{x}, \mathbf{y}) = \arg \max_{\mathbf{y}'} \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}') + \text{cost}(\mathbf{y}', \mathbf{y})$$

- Efficient decoding is possible when the features factor locally:

$$\mathbf{g}(\mathbf{x}, \mathbf{y}) = \sum_p \mathbf{f}(\mathbf{x}, \text{part}_p(\mathbf{y}))$$

- Efficient cost-augmented decoding requires that the cost function break into parts the same way:

$$\text{cost}(\mathbf{y}', \mathbf{y}) = \sum_p \text{local_cost}(\text{part}_p(\mathbf{y}'), \mathbf{y})$$

An Exercise

- If the features are such that we can use the Viterbi algorithm for decoding, what are some cost functions we could include inside an efficient cost-augmented decoding algorithm that's a very small change to Viterbi?

Structured Hinge

- Three different lines of work all arrived at this idea, or something very close.
 - Max-margin Markov networks (Taskar, Guestrin, and Koller, 2003)
 - Structural support vector machines (Tsochantaridis, Joachims, Hoffman, and Altun, 2005)
 - Online passive-aggressive algorithms (Crammer, Keshet, Dekel, Shalev-Shwartz, and Singer, 2006)
- Important developments in optimization techniques since then!
 - I'll highlight what I think it's most useful to know.

Max-Margin Markov Networks

- Taskar et al. actually work through a *dual* version of the problem.
 - Primal and dual are both QPs; exponentially many constraints or variables, respectively.
- Key trick: *factored dual*.
 - Enables kernelized factors in the MN.
 - Actual algorithm is sequential minimal optimization (SMO) for SVMs, a coordinate ascent method (Platt, 1999).
- The paper includes a generalization bound that is argued to improve over the Collins perceptron.
- Experiments: handwriting recognition, text classification for hyperlinked documents.

I'm Taking Liberties

- The M3N view of the world really thinks about outputs as configurations in a Markov network.
- They assume y corresponds to a set of random variables, each of which gets a label in a finite set.
- Their cost function is Hamming cost: “how many r.v.s do I predict incorrectly?”
 - This is convenient and makes sense for their applications. But it's not as general as it could be.

Structural SVM

- Tsochantaridis et al. (2005) – extends their 2004 paper.
- Slightly different version of the loss function:

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^N \xi_i$$

$$\text{s.t. } \forall i, \forall \mathbf{y}, \mathbf{w}^\top \mathbf{g}(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{w}^\top \mathbf{g}(\mathbf{x}_i, \mathbf{y}) \geq +1 - \frac{\xi_i}{\text{cost}(\mathbf{y}, \mathbf{y}_i)}$$

- Alternative version of cost-augmented decoding (“slack rescaling” as opposed to Taskar et al.’s “margin rescaling”)

Optimization Algorithms for SSVMs

- Taskar et al. (2003): SMO based on factored dual
- Bartlett et al. (2004) and Collins et al. (2008):
exponentiated gradient
- Tsochantaridis et al. (2005): cutting plane (based on dual)

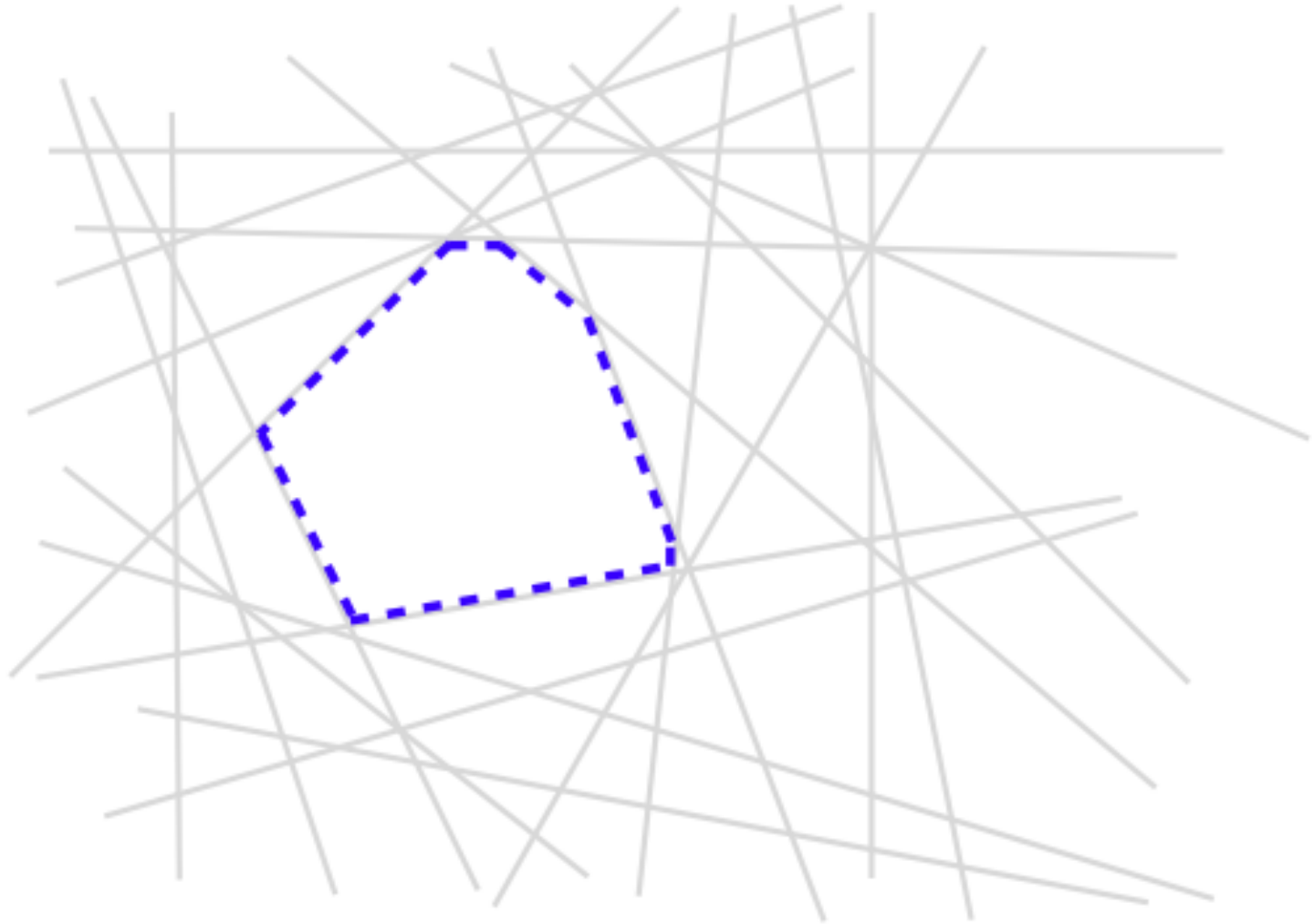
Cutting Plane

- There are exponentially many constraints!

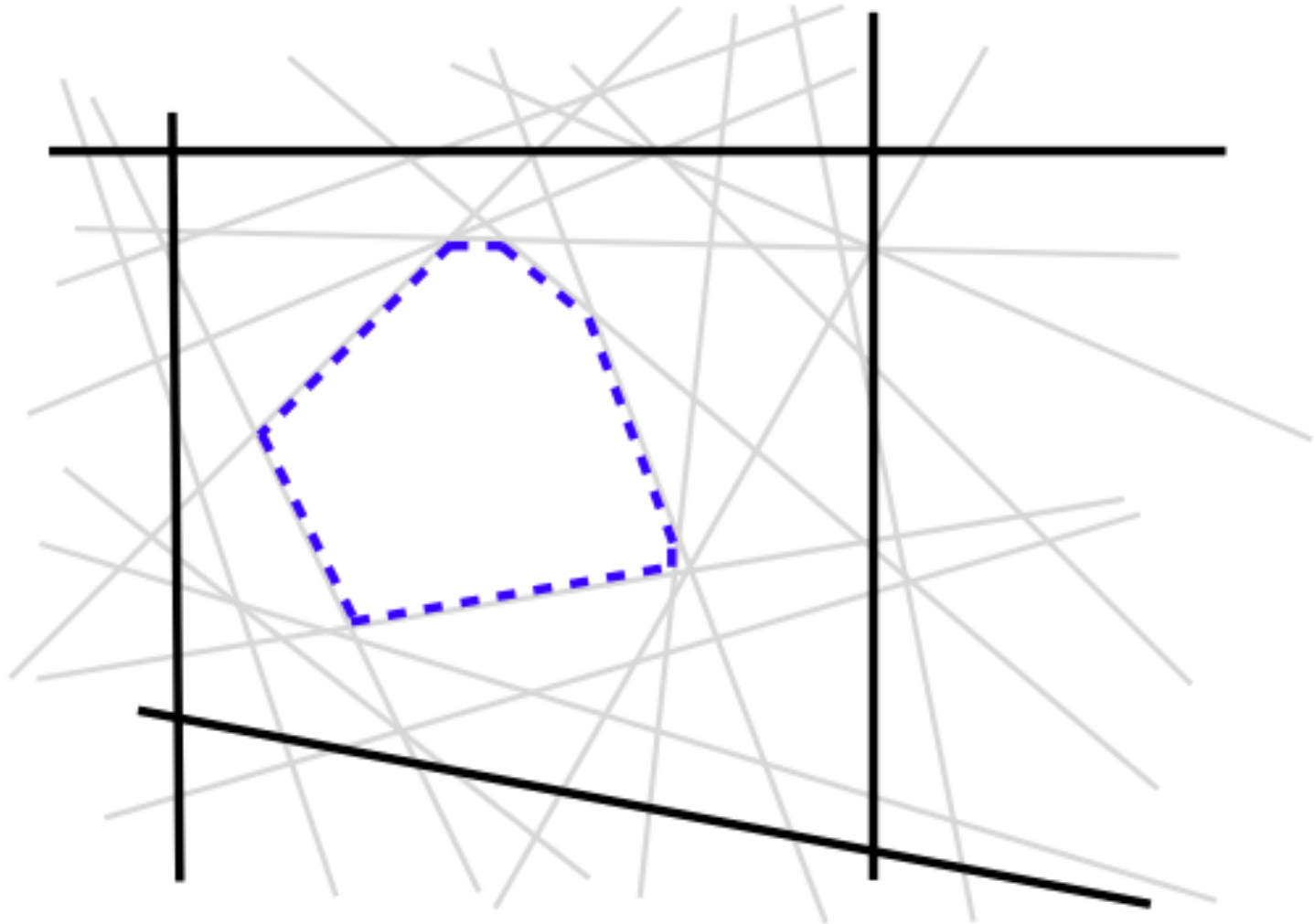
$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & \forall i, \forall \mathbf{y}, \quad \mathbf{w}^\top \mathbf{g}(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{w}^\top \mathbf{g}(\mathbf{x}_i, \mathbf{y}) \geq \text{cost}(\mathbf{y}, \mathbf{y}_i) - \xi_i \end{aligned}$$

- Instead of enumerating them all, let's dynamically add constraints as needed
- Iterate: solve a relaxation with a subset of constraints, then add most violated constraint

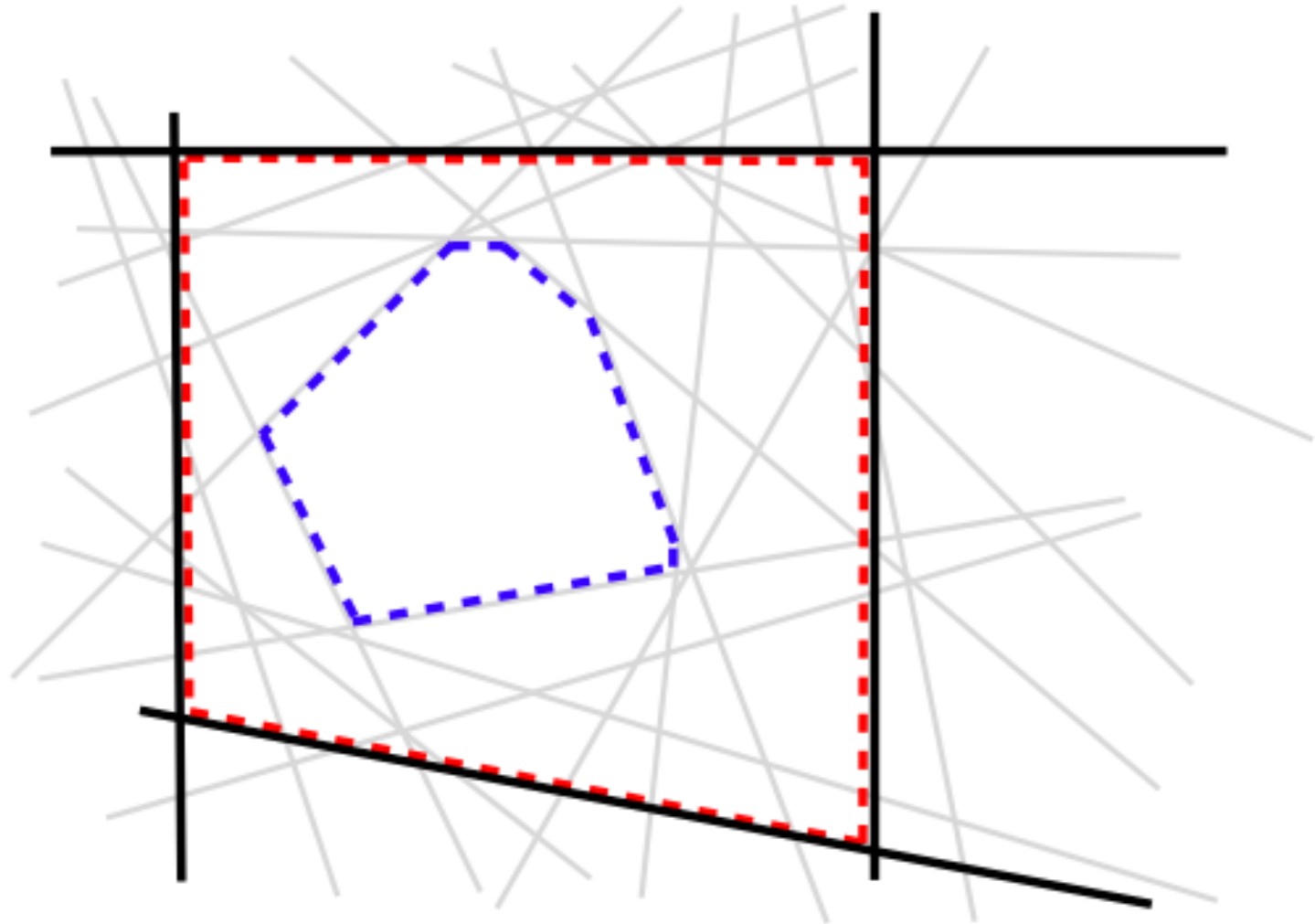
Cutting Plane



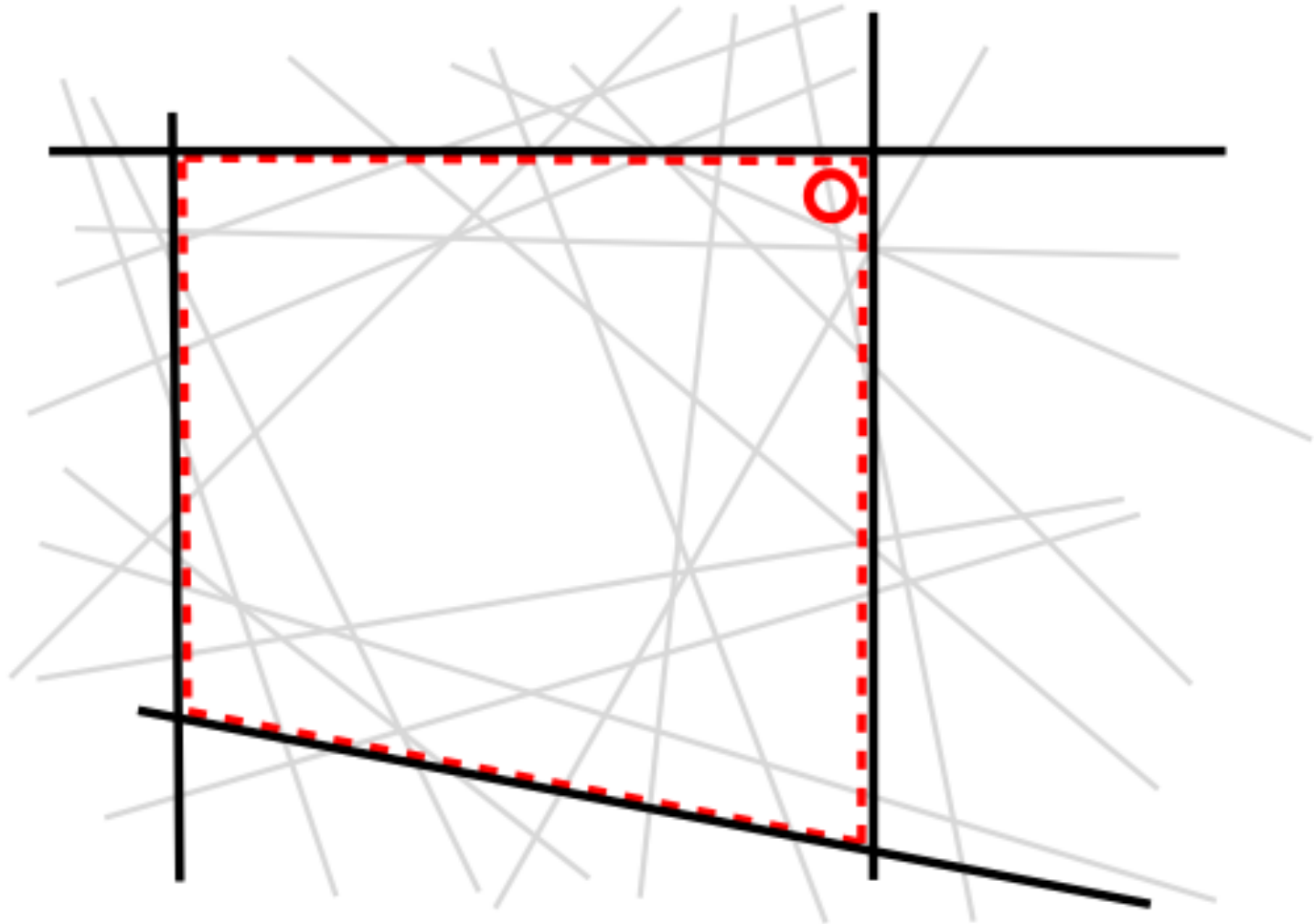
Cutting Plane



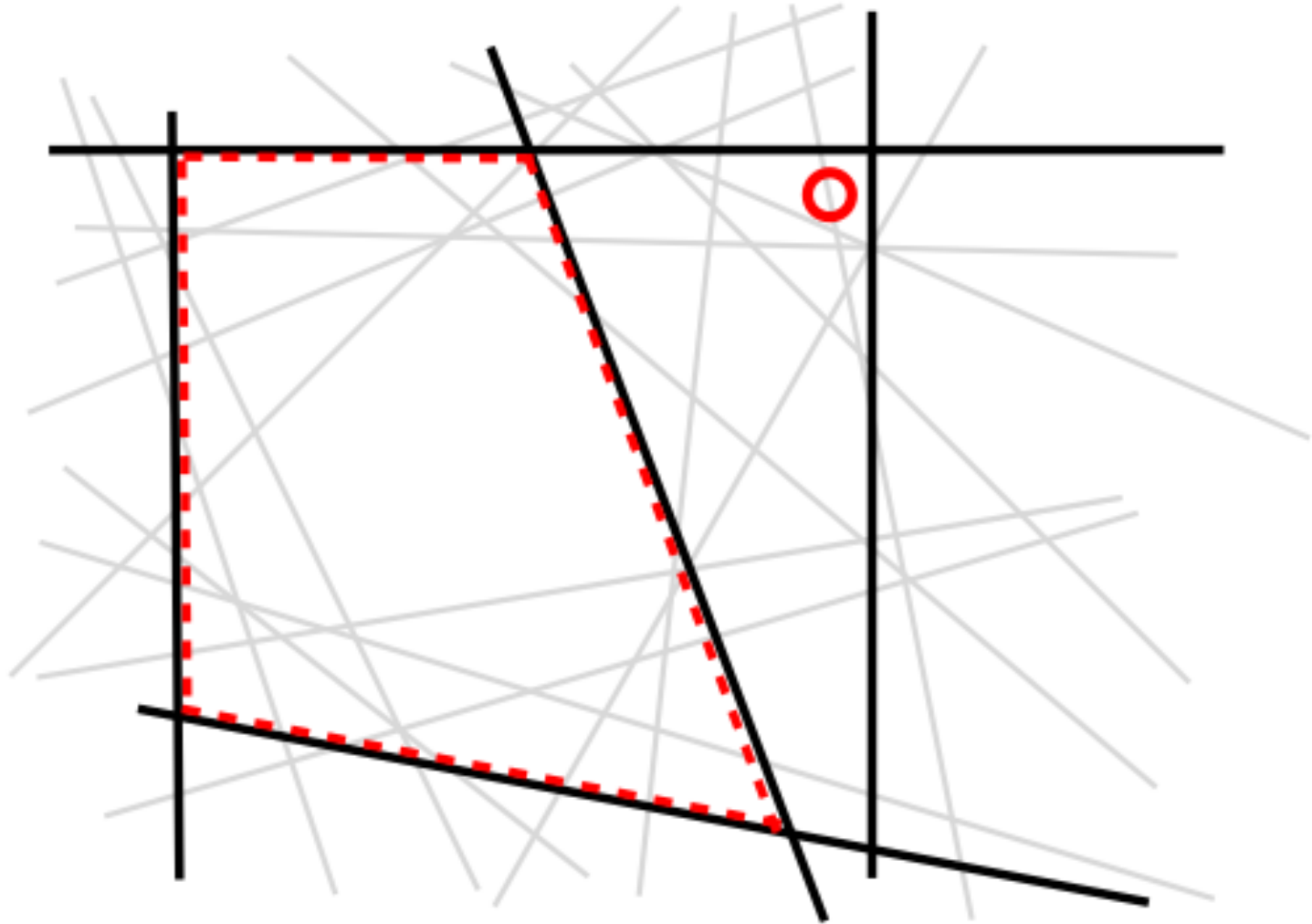
Cutting Plane



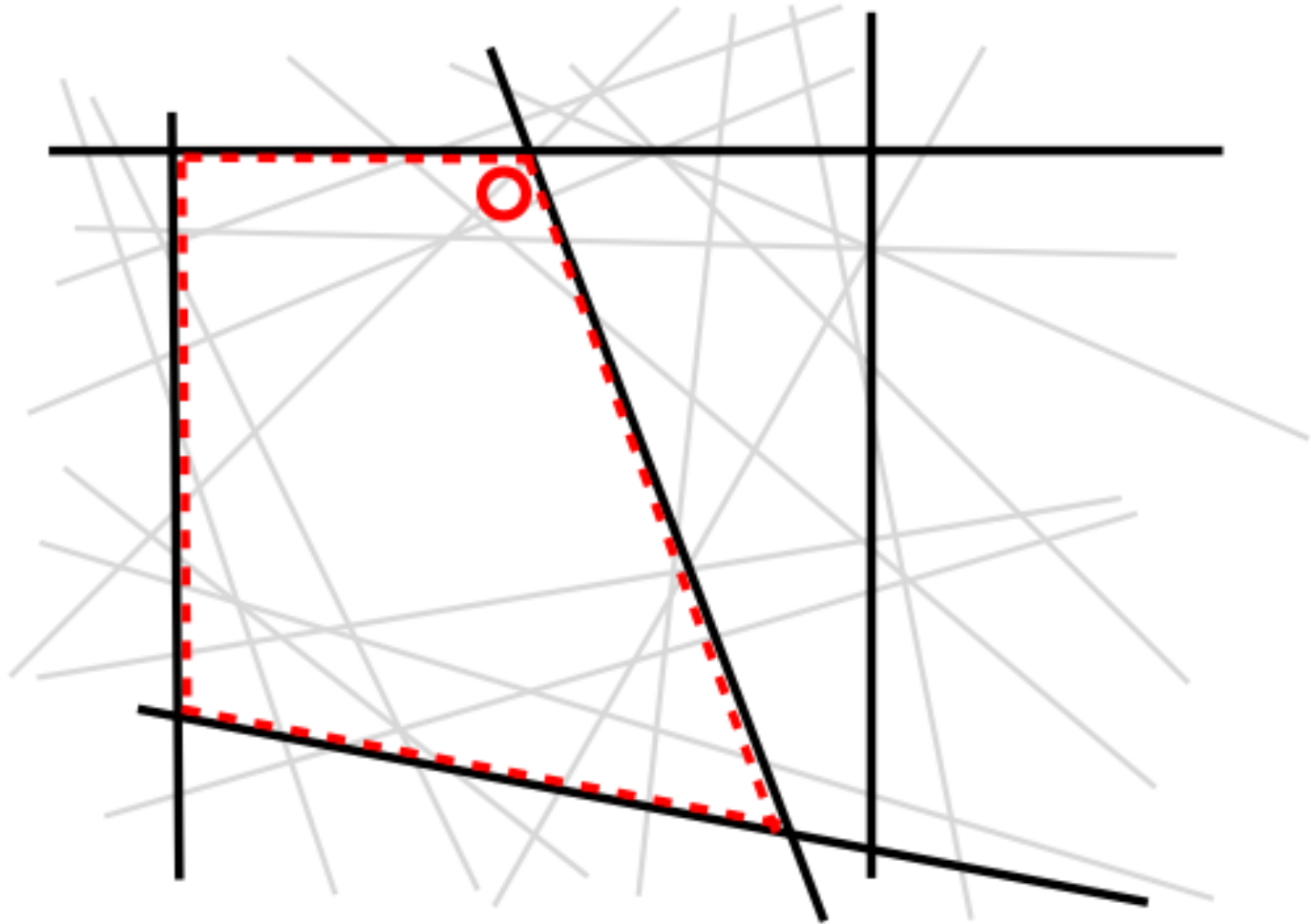
Cutting Plane



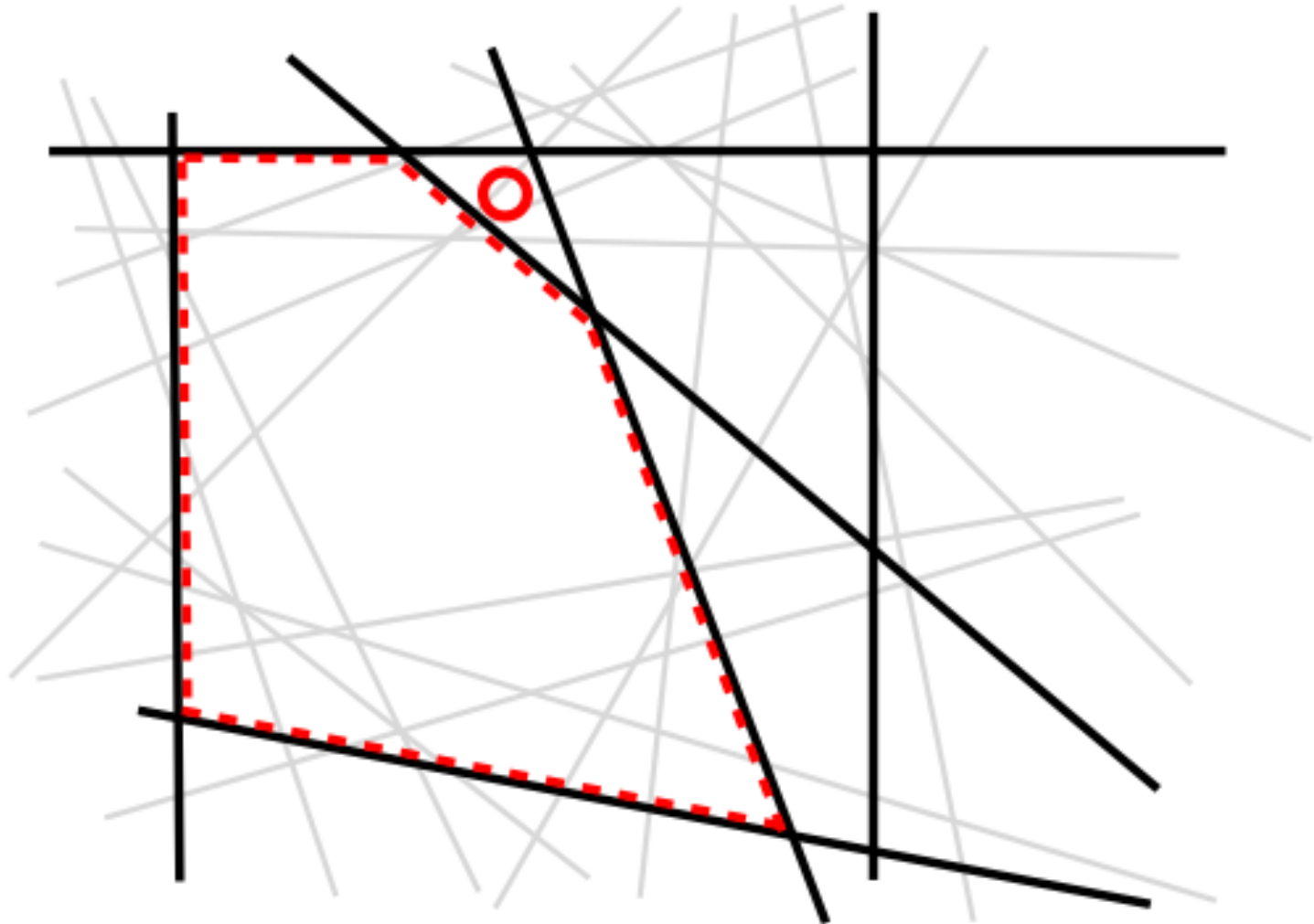
Cutting Plane



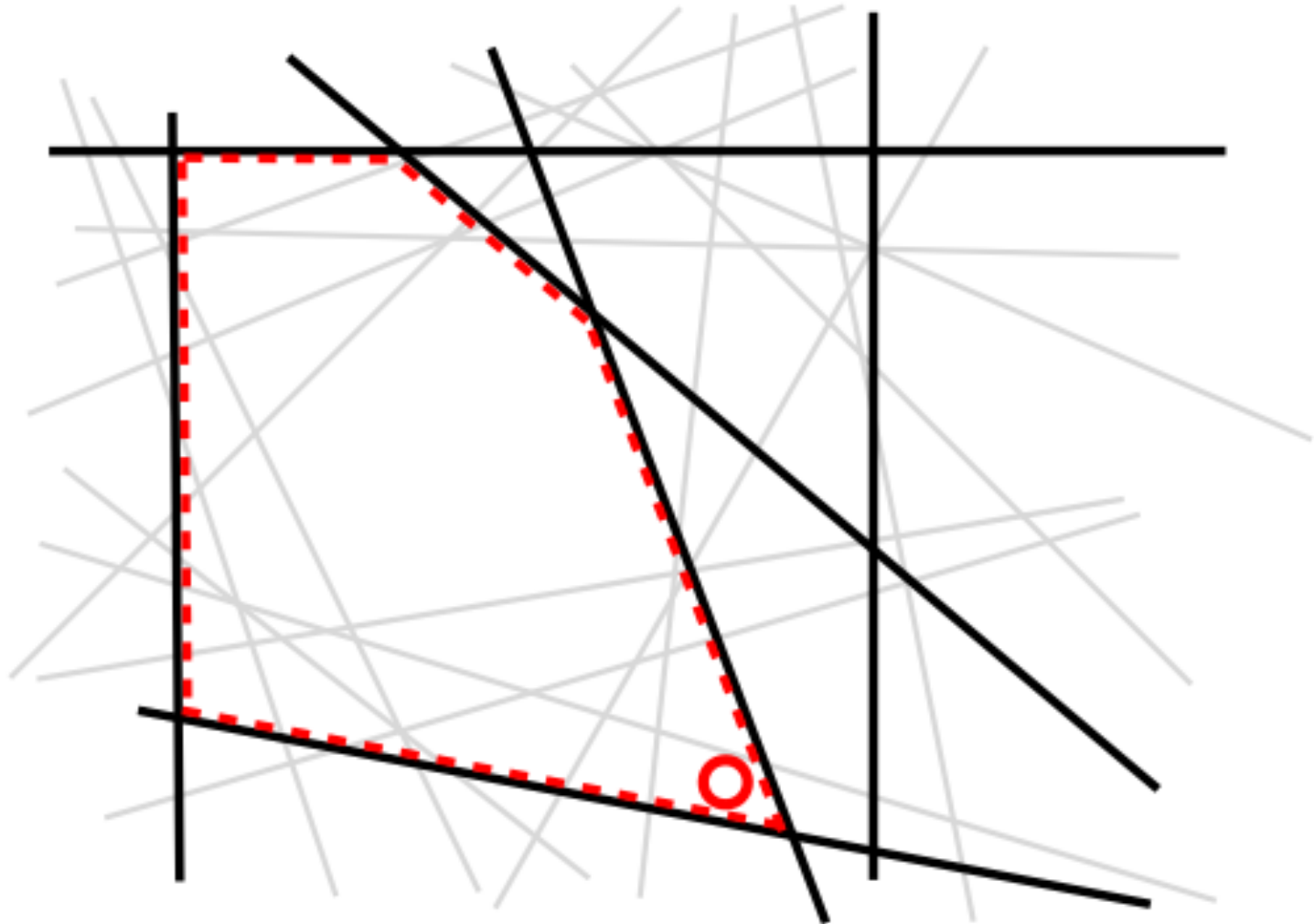
Cutting Plane



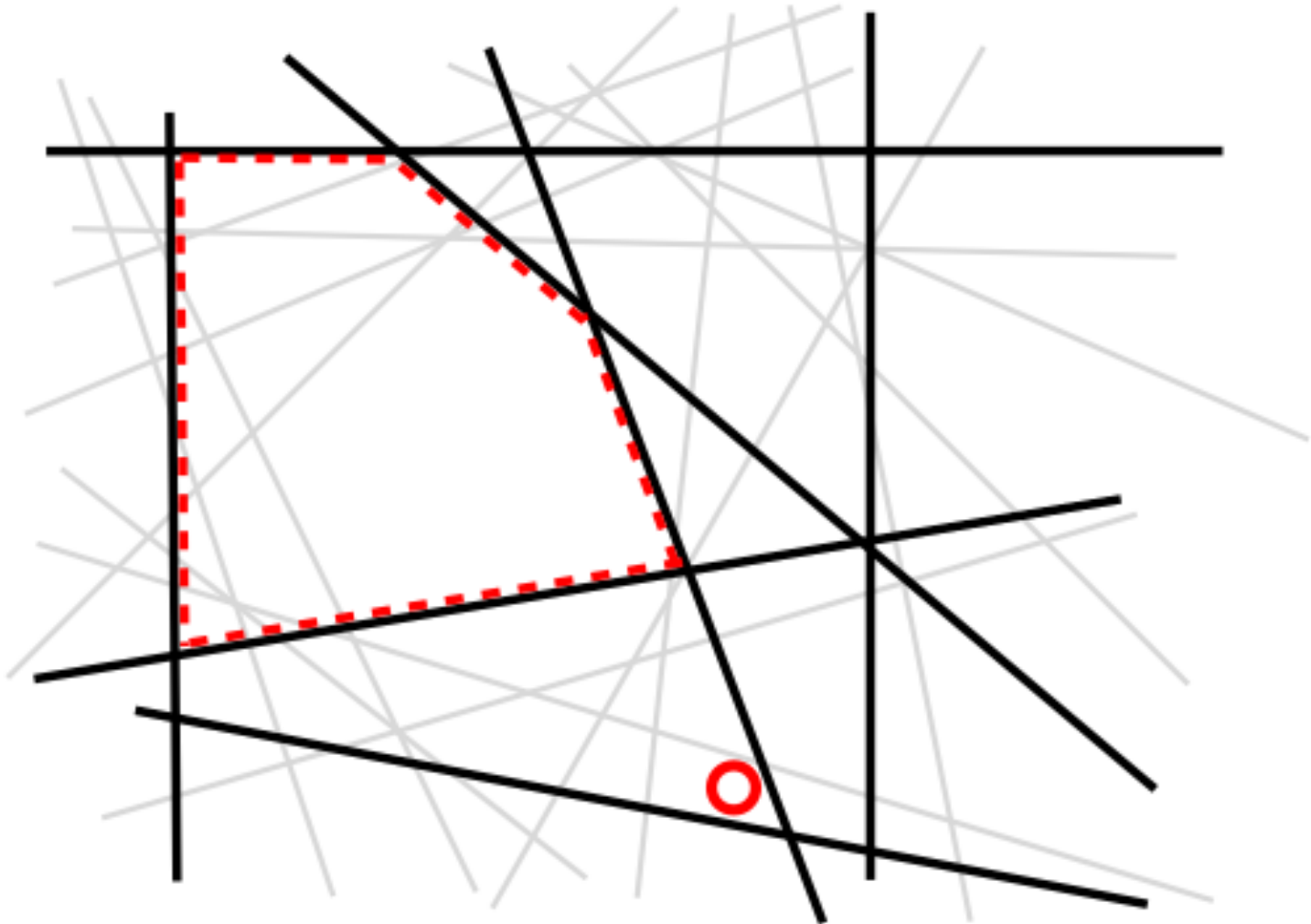
Cutting Plane



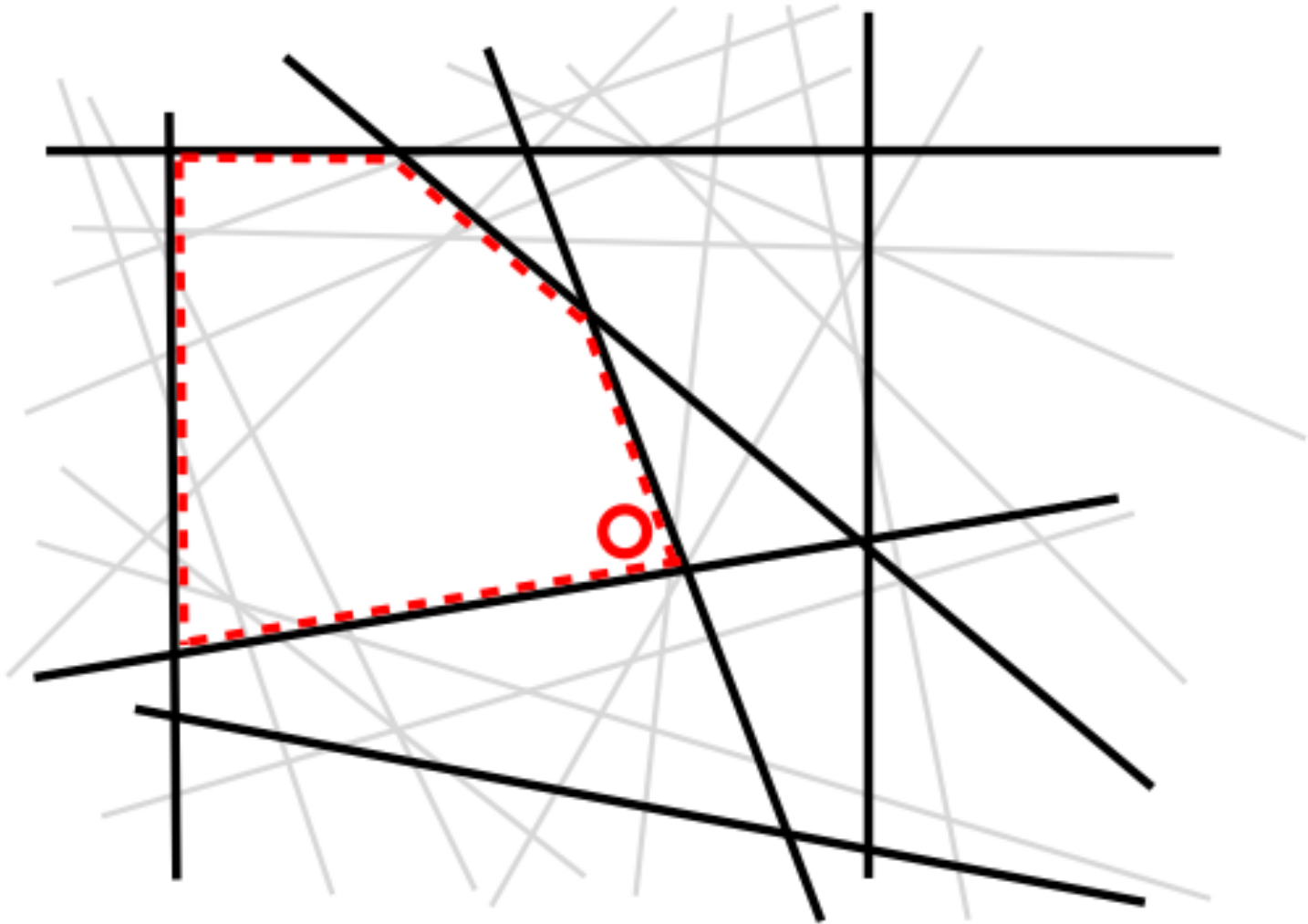
Cutting Plane



Cutting Plane



Cutting Plane



Optimization Algorithms for SSVMs

- Taskar et al. (2003): SMO based on factored dual
- Bartlett et al. (2004) and Collins et al. (2008): exponentiated gradient
- Tsochantaridis et al. (2005): cutting plane (based on dual)
- Taskar et al. (2005): dual extragradient

Easiest to use, in my opinion:

- Ratliff et al. (2006): (stochastic) subgradient descent
- Crammer et al. (2006): online “passive-aggressive” algorithms

Stochastic Subgradient Descent

- Unconstrained primal objective:

$$L(\mathbf{w}, \mathbf{x}, \mathbf{y}) = -\mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}) + \max_{\mathbf{y}'} \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}') + \text{cost}(\mathbf{y}', \mathbf{y})$$

- Resulting subgradient:

$$-g_j(\mathbf{x}, \mathbf{y}) + g_j(\mathbf{x}, \text{cost_augmented_decode}(\mathbf{w}, \mathbf{x}))$$

- Only requires loss-augmented decode! No need for marginals / summing (cf. CRF)

“Passive Aggressive” Learners

- Starting point is the perceptron, and the focus is on the step size.
- In NLP, people often use a specific instance called “1-best MIRA” (margin infused relaxation algorithm).
 - Sometimes with regular decoding, sometimes cost-augmented decoding.
- I do not understand the name (kind of I do)

Passive-Aggressive Update in a Nutshell (“1-best MIRA”)

- Given x (and y), perform decoding (or cost-augmented decoding) to obtain y' .
- To get the updated weights, solve:

$$\min_{\mathbf{w}'} \|\mathbf{w}' - \mathbf{w}\|_2^2$$

$$\text{s.t.} \quad \mathbf{w}^\top \mathbf{g}(x, y) - \mathbf{w}^\top \mathbf{g}(x, y') \geq \text{cost}(y', y)$$

- Closed form solution!
 - Essentially, a subgradient update with a closed-form step size.

Perceptron and PA

- The PA papers (e.g., Crammer et al., 2006) take a procedural view of online learning and prove convergence and regret-style bounds.
- An alternative view, described by Martins et al. (2010), derives the same updates via dual coordinate ascent.
 - Confusing name: it doesn't work in the dual!
 - More general: applies to many other loss functions, so you can get a closed-form step size for perceptron and CRFs.
 - Assumes L2 regularization; role of regularization constant C is very clear in the form of the update.

Dual Coordinate Ascent Update

$$\mathbf{w} \leftarrow \mathbf{w} - \underbrace{\min \left\{ \frac{1}{C}, \frac{L(\mathbf{w}, \mathbf{x}, \mathbf{y})}{\|\nabla_{\mathbf{w}} L(\mathbf{w}, \mathbf{x}, \mathbf{y})\|_2^2} \right\}}_{\text{step size}} \underbrace{\nabla_{\mathbf{w}} L(\mathbf{w}, \mathbf{x}, \mathbf{y})}_{\text{subgradient}}$$

- Assumes L2 regularization.
- 1-best MIRA is a special case with structured hinge loss.
- Can get regularization into perceptron this way (use perceptron loss).
- Can get closed-form step size for CRF stochastic GD.

Hinge Loss and Log Loss

- Hinge loss (M3N):

$$-\mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}) + \max_{\mathbf{y}'} \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}') + \text{cost}(\mathbf{y}', \mathbf{y})$$

- Log loss (CRF):

$$-\mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}) + \log \sum_{\mathbf{y}'} \exp \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}')$$

Aside: Probabilities *and* Cost?

- Hinge loss (M3N):

$$-\mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}) + \max_{\mathbf{y}'} \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}') + \text{cost}(\mathbf{y}', \mathbf{y})$$

- Log loss (CRF):

$$-\mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}) + \log \sum_{\mathbf{y}'} \exp \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}')$$

- “Softmax margin” (Gimpel and Smith, 2010):

$$-\mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}) + \log \sum_{\mathbf{y}'} \exp (\mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}') + \text{cost}(\mathbf{y}', \mathbf{y}))$$

Loss Functions You Know

Name	Expression of $L(\mathbf{w}, \mathbf{x}, \mathbf{y})$
Log loss (joint)	$-\log p(\mathbf{x}, \mathbf{y} \mid \mathbf{w})$
Log loss (conditional)	$-\log p(\mathbf{y} \mid \mathbf{x}, \mathbf{w})$
Cost	$\text{cost}(\text{decode}(\mathbf{w}, \mathbf{x}), \mathbf{y})$
Expected cost, a.k.a. “risk”	$\mathbb{E}_{p(\mathbf{Y} \mid \mathbf{x}, \mathbf{w})} [\text{cost}(\mathbf{Y}, \mathbf{y})]$
Perceptron loss	$\max_{\mathbf{y}'} \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}') - \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y})$
Hinge (margin rescaling version)	$\max_{\mathbf{y}'} \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}') + \text{cost}(\mathbf{y}', \mathbf{y}) - \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y})$

On Regularization

- In principle, this choice is orthogonal to the loss function.
- L2 is the most common starting place.
- L1 and other sparsity-inducing regularizers have some nice properties, but they can make optimization more complicated

Does this matter?

Practical Advice

- Features still more important than the loss function.
 - But general, easy-to-implement algorithms are quite useful!
- Perceptron is easiest to implement.
- CRFs and SSVMs usually do better.
- If the cost function factors locally, I recommend using a hinge loss and stochastic subgradient descent.
- Tune the regularization constant.
 - Never on the test data.