

A primer on structured prediction

1 What is structured prediction?

Supervised machine learning typically aims to learn a mapping from some input space \mathcal{X} to some output space \mathcal{Y} . In binary classification \mathcal{Y} is just $\{0, 1\}$. In multi-class classification \mathcal{Y} is $S = \{0, \dots, k\}$. However, for many computer vision tasks, the output space can be fairly complex. In semantic segmentation, for an $h \times w$ image, with k classes, the output space is S^{hw} , which is huge. In principle, we could consider this as just hw independent problems with an output space of S (e.g., by using convolutional network features to classify each pixel), but this independence assumption may not be justified: some combinations are inherently a lot more likely than others. For example, adjacent pixels are likely to have the same label. Such output spaces which are (a) large and (b) cannot be broken down into independent sub-problems because of underlying dependencies are often called “structured output spaces” and the problem is one of “structured prediction”.

In some cases like pose estimation, the underlying dependencies may not be known before hand. In such cases, we must also learn these dependencies.

2 Energy-based models

What we want is to model the conditional probability $P(\mathbf{y}|\mathbf{x})$. Here $\mathbf{y} = (y_1, \dots, y_n)$ where each y_i is a random variable. As an example, for semantic segmentation, each y_i would be a pixel.

Under fairly general conditions, we can write $P(\mathbf{y}|\mathbf{x}) \propto e^{-E(\mathbf{y}|\mathbf{x};\theta)}$, where E is some energy function with parameters θ . Sometimes we will write the negative of the energy as a *score function* $S(\mathbf{y}|\mathbf{x};\theta)$.

2.1 Special case: Graphical models

In principle, the energy function and resulting probability distribution can be completely general. However, in many applications, we have a significant amount of domain knowledge about what this probability distribution should look like. In particular, we can often make assumptions about the *conditional independence* of different variables. For example, in a semantic segmentation task, we might posit that given the input image and the labels of all *adjacent* pixels, a pixel’s label is *independent* of the labels of all other pixels.

Such conditional independence can be specified through a graphical model. A graphical model is a *family* of probabilistic models $P(\mathbf{y}|\mathbf{x})$ that all satisfy a set of conditional independence assumptions. These assumptions are expressed in the form of a graph, where each y_i is a node, and a given random variable is assumed to be independent of all other random variables *given its neighbors*. In general, for a graphical model, the energy distribution $E(\mathbf{y}|\mathbf{x};\theta)$ takes the form $E(\mathbf{y}|\mathbf{x};\theta) = \sum_{c \in \mathcal{C}} \phi_c(\mathbf{y}_c|\mathbf{x})$ where \mathcal{C} is the set of cliques in this graph. The functions ϕ_c are called *potentials*.

2.2 Inference

Suppose we know a good energy function. How do we use it to solve the problem? We might decide to look for the most probable labeling:

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}) = \arg \min_{\mathbf{y}} E(\mathbf{y}|\mathbf{x};\theta) = \arg \max_{\mathbf{y}} S(\mathbf{y}|\mathbf{x};\theta) \quad (1)$$

In a typical classification setting, the scoring function is just a score for each class, and doing this inference amounts to simply enumerating all the classes, looking at the corresponding score and picking the one with the highest score. However, in structured prediction problems, the output space is huge and cannot be enumerated. Instead, we must consider this as an optimization problem, and look towards methods in combinatorial (for discrete output spaces) or continuous (for continuous output spaces) optimization.

Because we cannot enumerate all possible outputs and must rely on doing optimization, it follows that not all forms of the energy function will lead to tractable inference. In fact, it is often the case that one needs to severely restrict the kind of energy functions we are looking at to keep inference tractable. For example, if the output space is continuous, E needs to be convex.

Even when the optimization is tractable, it is usually an iterative process. Thus, one starts with an initialization $\mathbf{y}^{(0)}$ and then updates it at every step. Exactly what this initialization is and how it is updated depends on the particular algorithm.

Some special cases related to graphical models:

1. If the graphical model is a completely disconnected graph with no edges, every variable is conditionally independent of others. In this case, the optimization can be done independently for each variable.
2. If the graphical model is a tree, then the optimization can be done using dynamic programming.
3. If each y_i can only be either 0 or 1 (e.g., segmenting just one class vs background), and if the potentials ϕ give higher scores (thus lower energies) when the nodes have the same label, then the optimization can be reframed as a *min-cut/max-flow* problem.

2.3 Learning

Now we want to consider how we can identify a good energy function. Often, if the prior is well understood, we can define the potentials by hand. However, an alternative is to *learn* the parameters θ .

As in traditional machine learning, this means we want to define a *loss function*. We can consider two different loss functions:

Negative log likelihood As with classification, we can use the negative log likelihood of the true labeling. Since $P(\mathbf{y}|\mathbf{x}) \propto e^{-E(\mathbf{y}|\mathbf{x})}$, we have $P(\mathbf{y}|\mathbf{x}) = \frac{e^{-E(\mathbf{y}|\mathbf{x})}}{\sum_{\mathbf{y} \in \mathcal{Y}} e^{-E(\mathbf{y}|\mathbf{x})}}$. The negative log likelihood of the true label is then:

$$-\log P(\mathbf{y}^*|\mathbf{x}) = E(\mathbf{y}^*|\mathbf{x}) - \log \sum_{\mathbf{y} \in \mathcal{Y}} e^{-E(\mathbf{y}|\mathbf{x})} \quad (2)$$

The first term is often easy to compute, but the second term is difficult, since it requires summing over all $\mathbf{y} \in \mathcal{Y}$. By analogy with statistical physics, this sum is called the partition function. Computing the partition function is tractable for a few kinds of graphical models, such as trees.

Margin-rescaled hinge loss An alternative loss function is the *margin-rescaled hinge loss* [2]. Intuitively, we want the true label y^* to score higher than every other labeling. We can try to enforce this. However, the model might then reach a solution where the true labeling scores just infinitesimally better than other labelings. Thus, we want the score of the true labeling $S(\mathbf{y}^*|\mathbf{x})$ to be higher than the score of every possible labeling $S(\mathbf{y}|\mathbf{x})$ by a *margin*. Structured prediction tasks have an additional property, which is that not all incorrect outputs are equally bad. For example, getting one keypoint wrong in a pose estimation problem is much better than getting all of them wrong. Let us assume that there is some function $\Delta(\mathbf{y}, \mathbf{y}')$ which measures the magnitude of the difference between \mathbf{y} and \mathbf{y}' . Then we want that wrong outputs for which $\Delta(\mathbf{y}^*, \mathbf{y})$ is high should have really low scores. Thus we want that:

$$S(\mathbf{y}|\mathbf{x}; \theta) \leq S(\mathbf{y}^*|\mathbf{x}; \theta) - \Delta(\mathbf{y}^*, \mathbf{y}) \quad \forall y \quad (3)$$

$$\Rightarrow S(\mathbf{y}|\mathbf{x}; \theta) + \Delta(\mathbf{y}^*, \mathbf{y}) - S(\mathbf{y}^*|\mathbf{x}; \theta) \leq 0 \quad \forall y \quad (4)$$

$$(5)$$

This constraint must be true for all y , which is equivalent to saying that it must be true for the y with the maximum left hand side:

$$\max_{\mathbf{y}} S(\mathbf{y}|\mathbf{x}; \theta) + \Delta(\mathbf{y}^*, \mathbf{y}) - S(\mathbf{y}^*|\mathbf{x}; \theta) \leq 0 \quad (6)$$

We can convert this equation into a loss:

$$L(\theta, \mathbf{x}, \mathbf{y}^*) = \max(0, \max_{\mathbf{y}} S(\mathbf{y}|\mathbf{x}; \theta) + \Delta(\mathbf{y}^*, \mathbf{y}) - S(\mathbf{y}^*|\mathbf{x}; \theta)) \quad (7)$$

This is the margin-rescaled hinge loss.

During training, we can sample images x and take a step along the gradient of $L(\theta, \mathbf{x}, \mathbf{y}^*)$. But there is a problem: the loss function involves a maximization over y , which as discussed above can be difficult because of the size of the output space. A careful look at the loss function reveals that what we actually need to maximize is the following: $S(\mathbf{y}|\mathbf{x};\theta) + \Delta(\mathbf{y}^*, \mathbf{y})$. This maximization looks suspiciously like the inference problem, barring the second term. If Δ is “nice” (e.g., it considers each pixel independently), then as long as we can do inference (Eq. (1)), we can do this maximization. In fact, this maximization is called “loss-augmented inference”.

3 Iterated refinement

A key issue in using energy-based models is that to keep inference (and thus learning) tractable, we have to make fairly strong assumptions about the energy function. Alternatively, we have to give up hope of exact inference.

To search for an alternative, some prior work argues that in all these models, the inference is typically an iterative process, and the energy function shows up only during this iteration. For example, if the inference procedure is gradient descent, then each step in the iteration amounts to:

$$\mathbf{y}^{(t)} \leftarrow \mathbf{y}^{(t-1)} - \nabla_{\mathbf{y}} E(\mathbf{y}|\mathbf{x};\theta) \Big|_{\mathbf{y}=\mathbf{y}^{(t-1)}} \quad (8)$$

. Instead of first learning an energy function and then doing approximate iterated inference as above, why not directly learn the iteration. For example, we can learn a function f that does this refinement.

$$\mathbf{y}^{(t)} \leftarrow f(\mathbf{y}^{(t-1)}, \mathbf{x}; \theta) \quad (9)$$

. This makes things easier because it is a straightforward feed-forward model that takes x and $y^{(t-1)}$ and can use them as it pleases.

Learning f is easy and can be done stage-wise. We only need three tuples of the form $(\mathbf{x}, \mathbf{y}^{(t-1)}, \mathbf{y}^*)$, where the first two are inputs and the last is the target. Let’s assume $\mathbf{y}^{(0)}$ is some default labeling. This allows us to create an initial dataset to train f . We then use f on a new set of labeled images to get $\mathbf{y}^{(1)}$, and create additional training tuples $(\mathbf{x}, \mathbf{y}^{(1)}, \mathbf{y}^*)$. We train f on the combined dataset. We keep repeating this as much as we want.

We can also train a different f for each time step. This gives us autocontext [3] instead of inference machines [1].

The disadvantage is that we have now lost the connection to probabilistic models.

References

- [1] D. Munoz. *Inference Machines: Parsing Scenes via Iterated Predictions*. PhD thesis, June 2013.

- [2] I. Tsochantaridis, T. Hofmann, T. Joachims, and Y. Altun. Support vector machine learning for interdependent and structured output spaces. In *ICML*, page 104, 2004.
- [3] Z. Tu. Auto-context and its application to high-level vision tasks. In *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*, pages 1–8. IEEE, 2008.