# Incorporating Human Domain Knowledge into Large Scale Cost Function Learning

**Markus Wulfmeier**
Oxford Robotics Institute
Department of Engineering Science
University of Oxford
markus@robots.ox.ac.uk

**Dushyant Rao**
Oxford Robotics Institute
Department of Engineering Science
University of Oxford
dushyant@robots.ox.ac.uk

**Ingmar Posner**
Oxford Robotics Institute
Department of Engineering Science
University of Oxford
ingmar@robots.ox.ac.uk

## Abstract

Recent advances have shown the capability of Fully Convolutional Neural Networks (FCN) to model cost functions for motion planning in the context of learning driving preferences purely based on demonstration data from human drivers. While pure learning from demonstrations in the framework of Inverse Reinforcement Learning (IRL) is a promising approach, we can benefit from well informed human priors and incorporate them into the learning process. Our work achieves this by pretraining a model to regress to a manual cost function and refining it based on Maximum Entropy Deep Inverse Reinforcement Learning. When injecting prior knowledge as pretraining for the network, we achieve higher robustness, more visually distinct obstacle boundaries, and the ability to capture instances of obstacles that elude models that purely learn from demonstration data. Furthermore, by exploiting these human priors, the resulting model can more accurately handle corner cases that are scarcely seen in the demonstration data, such as stairs, slopes, and underpasses.

## 1 Introduction

Manual handcrafting of cost functions for motion planning systems is an inherently complex and time consuming task. It requires high competency in the target area and expert knowledge about robotics systems and the applied algorithms. Ideally, robotic behaviour can be defined by untrained personnel, enabling task adaptation without involving highly trained experts. Inverse Reinforcement Learning (IRL) targets this problem by learning direct reward models from demonstration samples, and has been successfully applied to problems in a wide range of areas [1, 2, 3].

Recent advances exploit the ease of generating samples for learning from demonstration, and combining with with high capacity representations through Neural Networks in domains such as games [4] and autonomous driving [1]. While the possibility to create large amounts of training data without manual labelling enabled training large networks, corner cases still represent a challenge for deep learning as less training data is present. Our work targets initialising neural networks by employing human priors to improve performance in these cases.
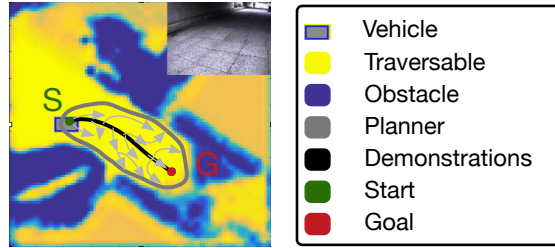
Figure 1: Illustration of sparse feedback, showing a demonstration trajectory on the spatial cost map around the vehicle, as well as the region explored by the planning algorithm. Error feedback is only created for the area surrounding sample trajectories.

While CNNs can be integrated straightforwardly into different domains, the IRL framework introduces additional challenges. One of the dominant influences is spatially sparse feedback from our objective function as displayed in Figure 1. For example, when training a model for image segmentation, the objective creates feedback for each individual pixel. When training with a loss based on Deep Inverse Reinforcement Learning on the other hand, error terms will focus on the region around the demonstration trajectories. These error terms are based on states visited by the demonstration samples and the planning algorithm, which inherently focuses around sample data. Our work addresses the shortcomings by pretraining the network towards a dense human-provided prior, to learn richer feature representations for untraversed areas and increase the network's ability to generalise.

We will show quantitatively that regression based pretraining improves prediction performance as well as classification performance for traversable terrain. Furthermore, we qualitatively present the advantages in the context of corner cases of manual cost functions, where pretraining followed by IRL-based training is able to recover more accurate and safe cost representations.

## 2    Related Work

A major share of early work in IRL focuses on small-scale scenarios and benchmarks [5, 6]. However, with recent technological advances, IRL approaches have been applied to larger state and feature spaces in real life applications [2, 1]. In particular, these techniques harness the potential of deep neural networks, learning rich representations that are able to model the relationship between the state of the environment and the reward structure implied by demonstrated behaviour.

When working with large state spaces and end-to-end learning of reward functions, exploration of the state space becomes more important to enable learning rich feature representations and improve on generalisation. Such capabilities are particularly important when learning driving behaviours from human demonstration: since the demonstrations cannot cover every possible driving scenario, it is necessary to ensure the vehicle can handle the unseen or scarcely seen "corner cases".

One paradigm that has been commonly employed in the past for deep neural networks is that of pretraining. In previous work, networks are pretrained in a greedy layer-wise fashion [7], by training each layer as a single layer unsupervised model (such as an autoencoder or restricted Boltzmann machine [8]), and using the hidden activations as the input to the next layer of unsupervised training. This pretraining has the effect of initialising the weights of the deep neural network to a region of the parameter space that is good for unsupervised tasks [9], from which the entire model can be fine-tuned to the specific classification task. This can be seen as being conceptually similar to the notion of *inductive transfer*, in which a model or representation learned for one task can be utilised or adapted to another.

In this work, we draw inspiration from these ideas, by performing pretraining of the deep IRL model with a manual cost map. In contrast to the unsupervised layer-wise training approach, however, we pretrain the network as a *regressor* to predict the manual cost map as output. This has the effect of initialising the network weights to a region of the parameter space which can accurately represent the manual cost map. This affords us the ability to inject domain knowledge into the network as a 'human prior', which is a crucial capability for self-driving vehicles to ensure that environments and behaviours that are absent from the expert demonstration data are handled gracefully.

## 3 Methods

### 3.1 Maximum Entropy Deep Inverse Reinforcement Learning

The goal of IRL is to infer the reward structure that underlies certain behaviours. The process can be defined under a Markov Decision Process framework $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \mathcal{T}, \gamma, r\}$, where $\mathcal{S}$ is the state space, $\mathcal{A}$ is the set of possible actions, $\mathcal{T}$ denotes the state transition model, $\gamma \in (0, 1]$ is a discount factor that moderates the influence of future rewards, and $r : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ is a function specifying the reward structure. As $r$ is unknown, it must be inferred from a set of demonstrations $\mathcal{D} = \{\varsigma_1, \varsigma_2, \ldots, \varsigma_N\}$, each of which is a sequence of state-action pairs $\varsigma_i = \{(s_1, a_1), (s_2, a_2), \ldots, (s_K, a_K)\}$ representing a sample trajectory.

The IRL model commonly needs to overcome two problems when reasoning about reward structures: suboptimality of sample trajectories given the underlying reward; and reward ambiguity, where multiple rewards can explain the same behaviour.

The Maximum Entropy (MaxEnt) approach to IRL explicitly addresses these concerns, by representing the expert behaviour as a distribution over demonstrated trajectories and assuming that this distribution has maximal entropy. The agent's behaviour under the policy $\pi_{\mathcal{D}}(a \mid s)$ maximises the reward given the current model. Consequently, the probability of any trajectory $\varsigma$ between specified initial and final states is proportional to the exponential of the reward along the path:

$$P(\varsigma \mid r) = \prod_{i=1}^{K} \pi_D(a_i \mid s_i) \propto \exp \left\{ \sum_{i=1}^{K} r_{s_i, a_i} \right\}. \tag{1}$$

The training loss in Maximum Entropy Deep Inverse Reinforcement Learning (MEDIRL) [10] consists of a data term, which maximises the log likelihood of the demonstration trajectories given the parametrised reward function, and a regularisation term:

$$\mathcal{L}(\theta) = \log P(\mathcal{D}, \theta \mid r(\theta)) = \underbrace{\log P(\mathcal{D} \mid r(\theta))}_{\mathcal{L}_{\mathcal{D}}} + \underbrace{\log P(\theta)}_{\mathcal{L}_{\theta}}. \tag{2}$$

The data term in Equation 2 can be split using the chain rule, into the gradient of the objective with respect to the reward, and the gradient of the reward with respect to the network parameters:

$$\begin{aligned} \frac{\partial \mathcal{L}_{\mathcal{D}}}{\partial \theta} &= \frac{\partial \mathcal{L}_{\mathcal{D}}}{\partial r} \frac{\partial r}{\partial \theta} \\ &= \underbrace{(\mu_{\mathcal{D}} - \mathbb{E}[\mu])}_{\substack{State Visitation \\ Frequency Matching}} \quad \underbrace{\frac{\partial}{\partial \theta} r(\theta)}_{Backpropagation} \quad . \end{aligned} \tag{3}$$

Here, the first term computes the difference in state visitation frequencies - how often specific states are traversed - between the demonstrations ($\mu_{\mathcal{D}}$) and the model; while the second term backpropagates this error through the network. It becomes obvious in this formulation that the area and therefore the types of features learned based on these error terms focus around the demonstration trajectories.

### 3.2 Model Architecture

The CNN architecture used for this work is the multi-scale fully convolutional network (MS-FCN) proposed in [1]. The input to the network is a 2D laser occupancy grid, and the output is a cost map specifying the reward for each location in the grid. The architecture consists of $5 \times 5$, $3 \times 3$, and $1 \times 1$ convolutional layers, including nonlinearities, max pooling and upsampling layers, and a separate branch to learn intermediate representations at a different scale, which are concatenated in the higher layers. For more details, the reader is referred to [1].

### 3.3 Training

Each iteration of training under the MEDIRL framework is composed of the following steps. First, the reward function (cost map) is computed by performing a forward pass of the network, and the MDP is solved for this current reward estimate. Using the policy $\pi$ computed by the MDP, we can obtain

the expected state visitation frequencies $\mathbb{E}[\mu]$ and the MaxEnt loss and gradients from Equations 2 and 3. Finally, these errors are backpropagated through the network to obtain the weight updates.

In the pretraining setup, the CNN first learns to predict the manual cost map as output - which results in a dense map of gradients over the full area - and then fine-tuned with the above steps. This is discussed in further detail in the next section.
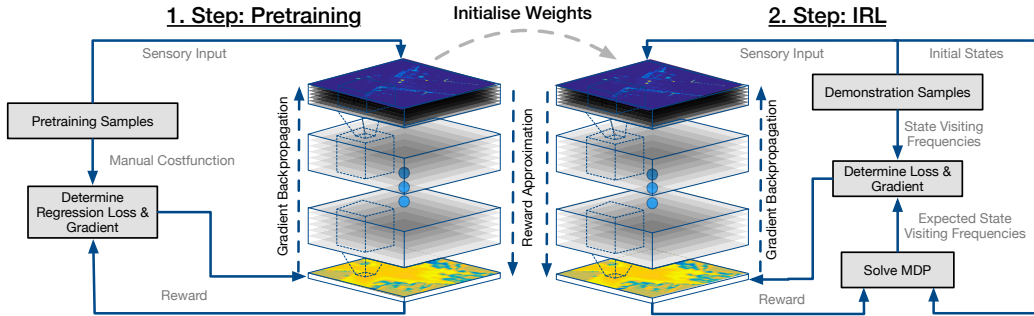
## 3.4 Incorporating Human Priors



Figure 2: Schema for additional network pretraining, where the model learns to regress to a manual prior cost map. Subsequently the network is fine-tuned to predict the reward under the MEDIRL framework.

Due to the iterative nature of planning and refining the cost model in the training process, the approach described in Section 3.1 tends to focus on discriminative features around highly visited states as depicted in Figure 1. Features for terrains that are neither explored by demonstration samples nor the planning step will never be formed by backpropagation of error terms through the network. Therefore, the model has to generalise based on similarity to more commonly traverse areas, resulting in artifacts and noisy reward maps. This situation is enhanced in scenarios where feature representations for similar places differ based on their position in state space, as is the case in the dataset used in this work [1]. The LIDAR scan points in this setup will be spread sparser at greater distance from the car, resulting in spatially variant representation.

In order to address this shortcoming, we suggest training the network as a regressor towards a prior cost map. These cost maps can be automatically generated from the laser input data based on manually handcrafted features, which enables us to utilise the availability of large amounts of data without human labelling efforts. However, the principal benefit is the ability to explore all features relevant to generating the cost function, leading to better generalisation in areas with greater distance from the demonstration trajectories.

## 4 Experiments

We evaluate all approaches on the large scale dataset presented in [1] and investigate how performance progresses between training from random initialisation and employing prior human knowledge.

## 4.1 Data

The dataset consists of 25,000 demonstration samples of urban driving from 13 different drivers. The data was collected in the inner city of Milton Keynes and includes driving around different types of obstacles, including but not restricted to: bollards, green patches, bike racks, slopes, cars, and underpasses. To increase the performance of our approaches as path planning cost maps, we added common preprocessing steps, such as normalising the input data, and trained on shorter trajectories that are more representative for the motion primitives employed in our motion planning framework.

## 4.2 Prediction Performance

To evaluate how well the trained models approximate human behaviour, we use two common metrics: the negative log-likelihood (NLL) and the modified Hausdorff Distance (MHD) [11]. The first metric is representing how likely the expert demonstrations are given the current cost function and the latter is a spatial metric for how close the demonstrations are to samples drawn from a probability distribution given the cost map. The learned models generally outperform the handcrafted approach as increasing the probability of the demonstration samples is inherently part of their objective function, in contrast to the manual cost function. As Table 1 displays pretraining improves our ability to predict where people are more likely to drive.

| Metric | NLL | MHD |
|---|---|---|
| Manual cost function | 56.402 | 0.286 |
| wo pretraining | 47.535 | 0.218 |
| w pretraining | **46.767** | **0.182** |

Table 1: Evaluation of cost functions for urban driving under the negative log-likelihood (NLL) and Modified Hausdorff Distance (MHD) metrics. Lower numbers represent models that are approximating human behaviour with higher precision.

## 4.3 Classification Performance

One major drawback for evaluating all approaches in a real world setup is the absence of absolute ground truth for the cost map. To overcome this impediment, we evaluate the approach based on its performance as a classifier for feasible example trajectories. Given a set of traversable trajectories - taken from our test set - and a set of artificial collision trajectories, we can analyse the accuracy of the model in a classification setup for traversable trajectories. In this setup, traversable terrain is taken as the positive class, meaning that a high precision model is conservative when assigning traversability.
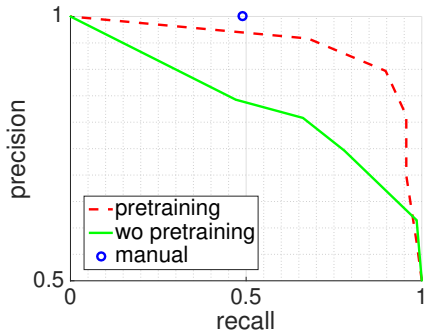


Figure 3: Precision Recall (PR) Curves for Trajectory Classification. The manual cost function has high precision but low recall, meaning that it is safe but conservative and will falsely classify a significant number of feasible trajectories as untraversable. Applying human priors in the pretraining step enables a significant gain in precision towards the baseline. This method approaches the precision of the manual cost function while strongly exceeding it in recall.

While utilising human priors already improves accuracy for prediction, its principal gain lies in being able to improve spatial generalisation and the ability to learn more robust cost functions. Figure 3 depicts the Precision-Recall curves for networks with and without pretraining. The handcrafted cost function does not include a threshold parameter and is therefore represented as a point. While this approach is manually designed to be conservative and enables us to operate at maximum precision, it falsely rejects much of the terrain as untraversable. The learned cost functions enable us to find possible paths in many situations when the manual cost function will get stuck. When introducing human prior knowledge into the training process, the approach achieves a significant gain in precision compared to random initialisation. Hence, utilising this knowledge is an important step towards robust application of learned cost maps.

| Scenario | Manual cost function | wo pretraining | w pretraining |
|---|---|---|---|
| Stairs | | | |
| Bollards | | | |
| Grass | | | |
| Underpass | | | |
| Slope | | | |

Table 2: Corner Cases for the Cost Function. The photos include views from the front facing camera module with the vehicle represented gray rectangle driving towards the right side of each cost map. Obstacles are represented in blue while yellow depicts traversable terrain.

## 4.4 Qualitative Assessment of Learned Cost Functions

Table 2 represents various situations emphasising shortcomings of the manual cost function. One point of high importance is, that all learned cost functions show additional obstacles starting at about 13 meters distance from the vehicle position, which is the length of the demonstration trajectories. By choosing the length of the trajectories, we define how far we trust our perception systems. Since features in distant areas are only traversed but planner and not demonstration samples in the training process, they will classified as untraversable with high probability.

The principal rules behind the handcrafted cost function are based on a threshold on the height range of detected points within a cell and the expansion of obstacles by the size of the vehicle to enable point based planning. This can lead to inaccuracies in the presence of slopes, which can exceed the threshold and will be shown as untraversable, and the same can occur for underpasses, where scans from ceiling and floor result in a high height range. Stairs on the other hand can still fit within the same threshold but present obstacles for any vehicle since they cannot be traversed due to their indiscontinuity. Bollards that are extended slightly too far will seem untraversable and areas such as grass might look very similar in features to pathways but should not be traversed.

While the randomly initialised network already learns to represent the main obstacles and traversable areas, it results in some noisy areas and artificial obstacles. Based on human prior domain knowledge, the network learns to refine the representation and is significantly more robust, learns to represent distinct obstacle boundaries and displays fewer artifacts. The approach learns that slopes are traversable, while stairs are not and extends obstacle boundaries only as far as necessary for safe traversal as seen in the respective cases in Table 2.

We conjecture that without pretraining, a lot of the expressive power of the network is employed to learn the very specific representation focused around the demonstration trajectories. If we instead pretrain the model to predict an existing cost map, the remainder of the learning process is able to generalise broader and better capture some of the corner cases described previously.

## 4.5 Implementation Details

When performing path planning, a threshold has to be determined to define untraversable and unsafe terrain. To simplify this step we normalise the output by sending it through a sigmoid. After first trials with Rectified Linear Units (ReLU) as activation functions within the network, the sigmoid saturates early on and training slows down, we exchanged all activation functions with sigmoid units. Furthermore, we introduce early stopping to increase generalisation performance; since following the pretraining step, the model requires fewer iterations to converge.

## 5 Conclusions and Ongoing Work

We develop an approach for incorporating human prior knowledge into cost function learning in the context of IRL. By utilising human priors, we can improve on prediction of demonstrator trajectories. Adding the pretraining step enables us to start from a significantly more desirable network initialisation; giving a richer representation for features at arbitrary position and enabling more effective generalisation. The approach improved classification performance for traversable terrain and results in more robust and distinct cost representations. Our evaluation furthermore depicts the advantages in specific corner cases of our environment given by slopes, bollards, underpasses and stairs.

Ongoing research looks into employing variations of Progressive Neural Networks [12], which represent another opportunity to exploit human priors by switching to modelling the error/residual between the human given prior and the true cost function. Different variations of the approach can focus on directly reusing the learned features of the trained model or simply focus on approximating the difference between the manual cost function and the optimal function to describe human behaviour.

The focus shifts in context of the latter from learning all representations needed to approximate a function to learning the difference between the manual cost function and the one best suited to explain human behaviour. This approach benefits from expert domain knowledge such that it only needs to embody the variations that the expert did not consider.

## 6 Acknowledgements

## References

[1] Markus Wulfmeier, Dominic Zeng Wang, and Ingmar Posner. Watch this: scalable cost-function learning for path planning in urban environments . In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2016. arxiv preprint: http://arxiv.org/abs/1607.02329.

[2] Chelsea Finn, Sergey Levine, and Pieter Abbeel. Guided cost learning: Deep inverse optimal control via policy optimization. *CoRR*, abs/1603.00448, 2016.

[3] Quoc Phong Nguyen, Bryan Kian Hsiang Low, and Patrick Jaillet. Inverse reinforcement learning with locally consistent reward functions. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 1747–1755. Curran Associates, Inc., 2015.

[4] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.

[5] Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, and Anind K Dey. Maximum entropy inverse reinforcement learning. In *AAAI*, pages 1433–1438, 2008.

[6] Andrew Y Ng and Stuart J Russell. Algorithms for inverse reinforcement learning. In *Proceedings of the Seventeenth International Conference on Machine Learning*, pages 663–670. Morgan Kaufmann Publishers Inc., 2000.

[7] Yoshua Bengio, Pascal Lamblin, Dan Popovici, Hugo Larochelle, et al. Greedy layer-wise training of deep networks. In *Advances in neural information processing systems*, 2007.

[8] Adam Coates, Honglak Lee, and Andrew Y Ng. An analysis of single-layer networks in unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*, pages 1–9, 2010.

[9] Dumitru Erhan, Yoshua Bengio, Aaron Courville, Pierre-Antoine Manzagol, Pascal Vincent, and Samy Bengio. Why does unsupervised pre-training help deep learning? *Journal of Machine Learning Research*, 11(Feb):625–660, 2010.

[10] Markus Wulfmeier, Peter Ondruska, and Ingmar Posner. Maximum entropy deep inverse reinforcement learning. In *Neural Information Processing Systems Conference, Deep Reinforcement Learning Workshop*, volume abs/1507.04888, Montreal, Canada, 2015.

[11] Kris M. Kitani, Brian D. Ziebart, James Andrew Bagnell, and Martial Hebert. Activity Forecasting. In Andrew Fitzgibbon, Svetlana Lazebnik, Pietro Perona, Yoichi Sato, and Cordelia Schmid, editors, *Computer Vision – ECCV 2012*, number 7575 in Lecture Notes in Computer Science, pages 201–214. Springer Berlin Heidelberg, 2012.

[12] Andrei A. Rusu, Neil C. Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *CoRR*, abs/1606.04671, 2016.